

Algorithms in Bioinformatics

De Novo Genome Assembly
and Short Read Mapping

Carlo M. Lapid

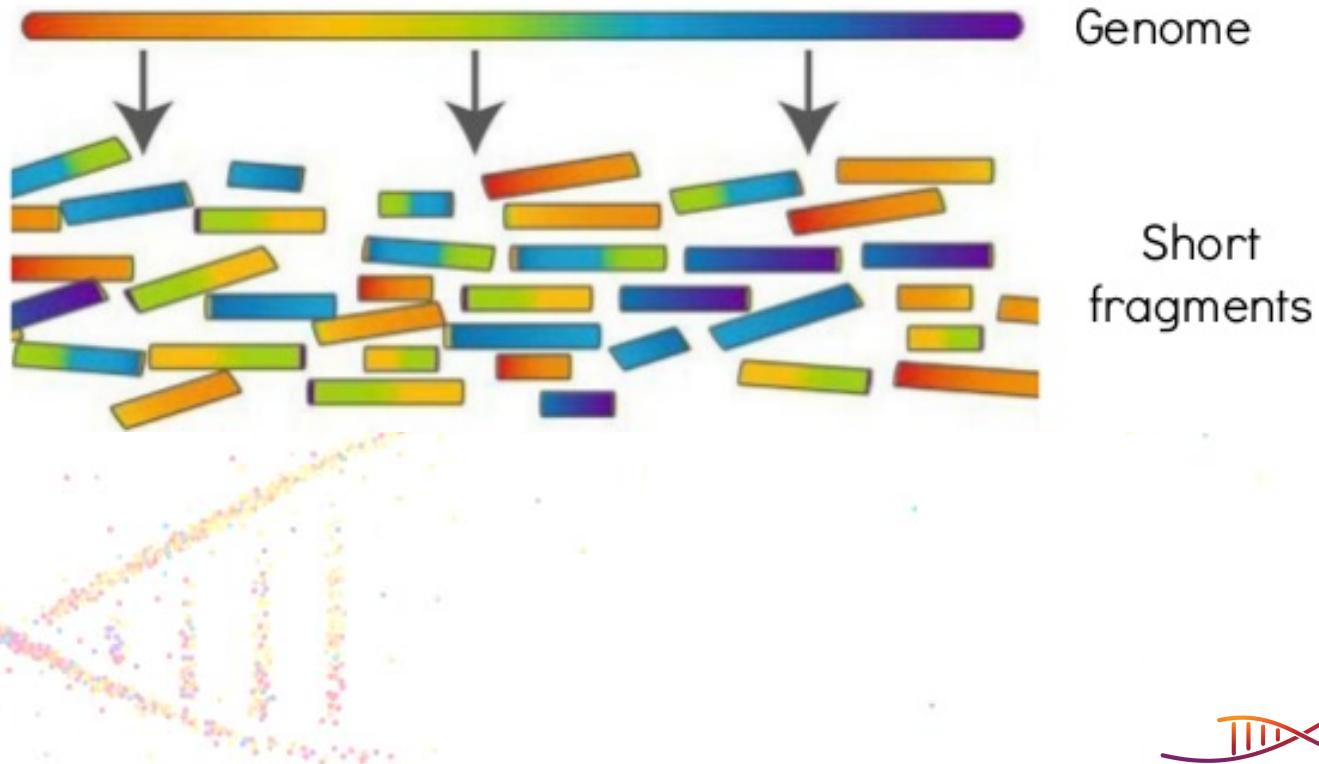
*Core Facility for Bioinformatics
Philippine Genome Center*

Algorithms

1. De novo genome assembly
2. Mapping short reads to a reference genome

What is assembly?

Next-generation sequencing creates millions of short, unorganized reads with no clear biological meaning or context.



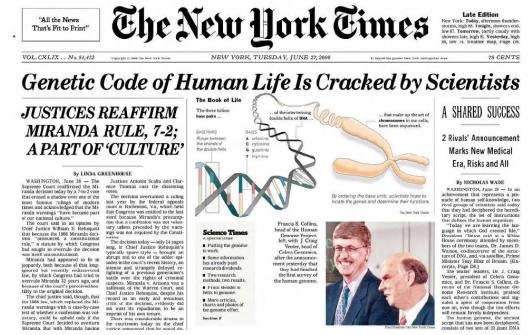
What is assembly?



Copies of genome



Sequencing reads



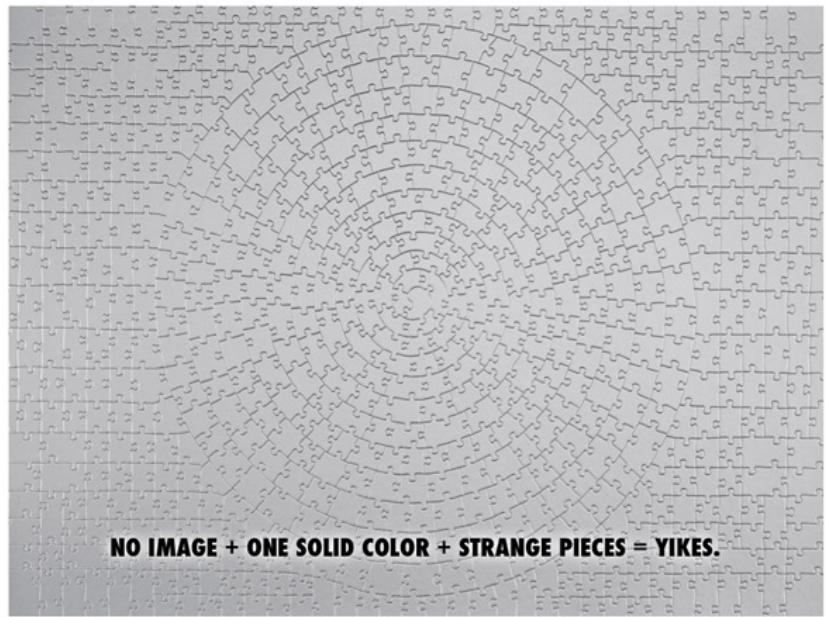
Assembled genome



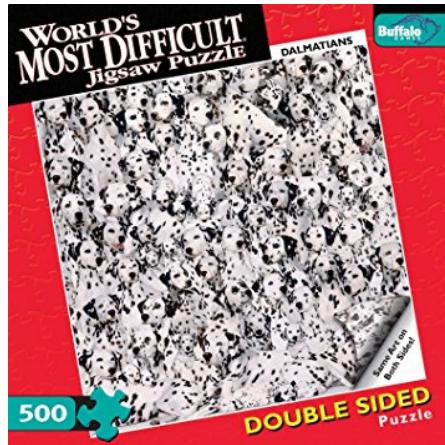
What is assembly?

- De novo genome assembly is like solving a jigsaw puzzle...
 - with millions of pieces,
 - multiple copies of each piece,
 - pieces that are missing,
 - pieces that fit together but shouldn't,
 - pieces that don't fit together but should
 - pieces that don't belong to the puzzle,
 - and no picture on the box.

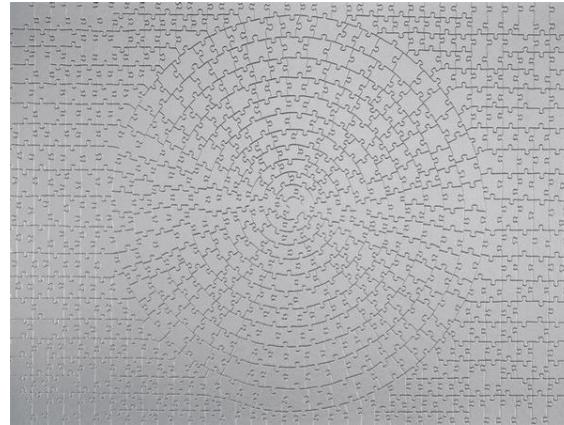
WANT A CHALLENGE?
YOU GOT IT.



What is assembly?



Reference-based



De novo

- Assembly can be performed by mapping reads to a known reference genome, or accomplished *de novo*, i.e. “from the beginning”.
- Why *de novo*?
 - Sometimes an appropriate reference genome isn't available
 - Inappropriate reference genome can mislead assembly
 - *De novo* techniques can uncover variation that might be absent in a reference genome

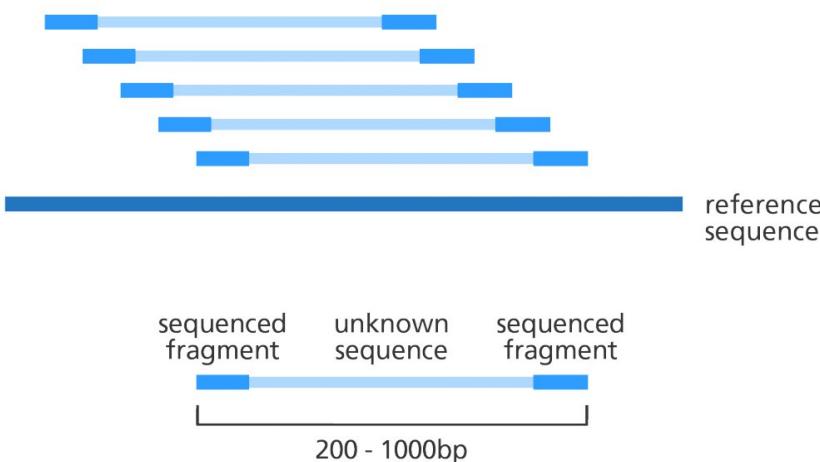
De novo assembly

Single-end reads



Depending on the sequencing protocol and technology, sequencing reads can be *single-end* or *paired-end*.

Paired-end reads

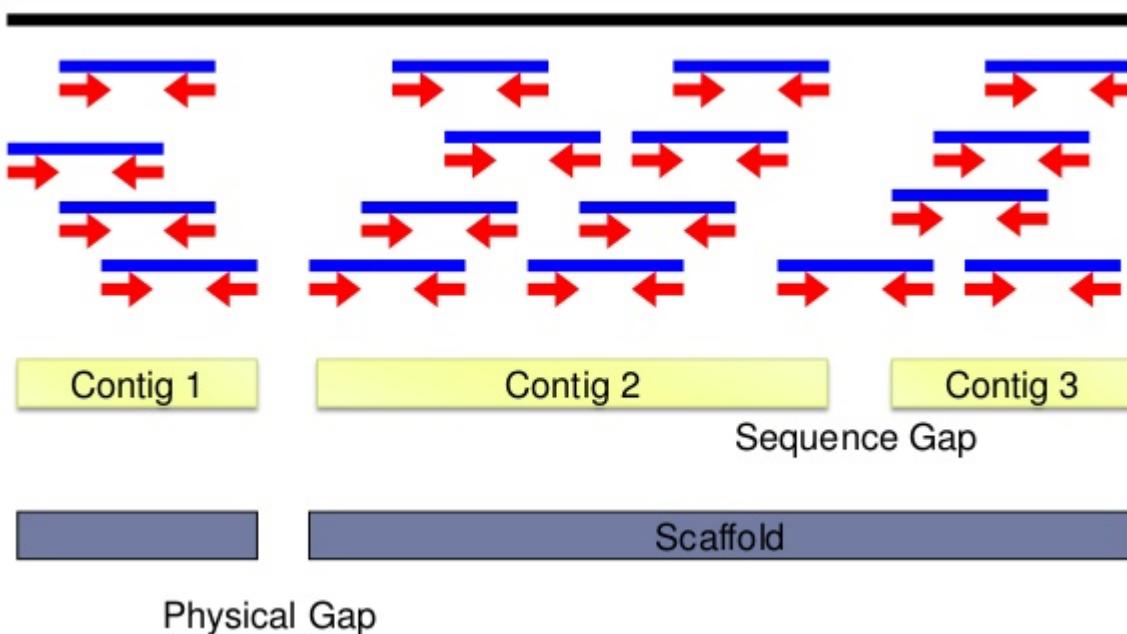


Paired-end reads are sequenced from opposite ends of the same DNA fragment, and the approximate genomic distance between them is known.

De novo assembly

Contigs are contiguous sequences assembled from overlapping reads.

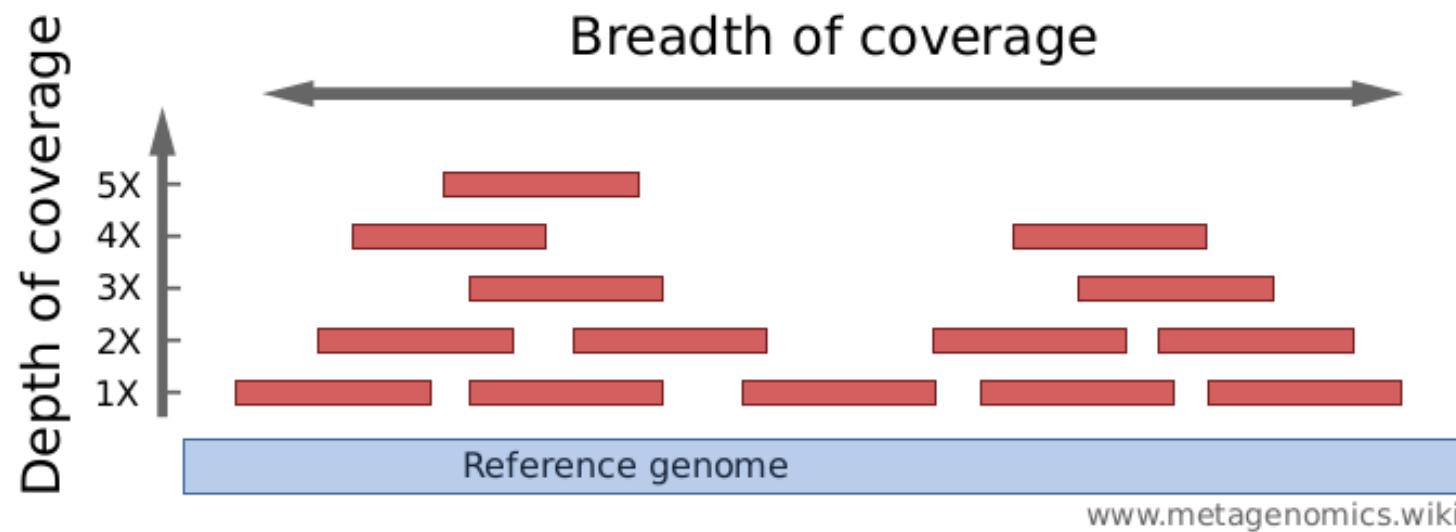
Scaffolds are assemblies of contigs, separated by unsequenced gaps of approximately known size.



De novo assembly

Depth of coverage refers to how often a base of a genome is sequenced, i.e. “covered” by reads.

Breadth of coverage is the percentage of bases of the target genome covered with a certain depth.



Assembly: an example

Exercise: Try to assemble these five reads into a contig.

Read 1: **A**C**G****G****A****T****A****C**

Read 2: **G****T****C****C****T****T****G****C**

Read 3: **A****T****A****C****A****G****G****A**

Read 4: **C****A****T****A****G****T****C****C**

Read 5: **A****G****G****A****C****A****T****A**

Assembly: an example

Read 1:

A C G G A T A C

Read 3:

A T A C A G G A

Read 5:

A G G A C A T A

Read 4:

C A T A G T C C

Read 2:

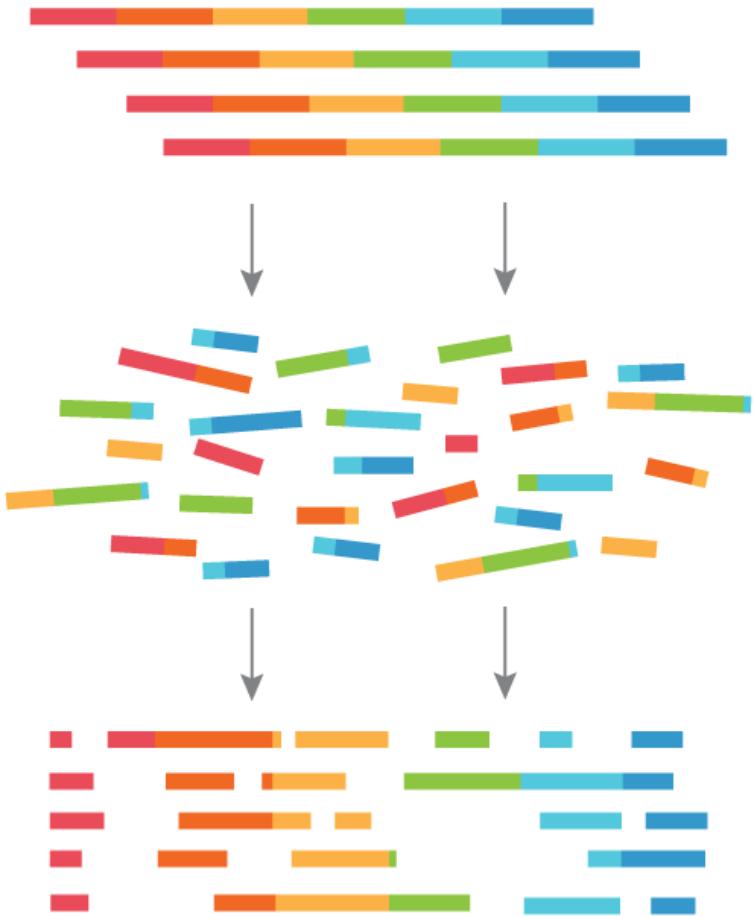
G T C C T T G C

Assembly:

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

What steps did you instinctively take to arrive at this solution?

De novo assembly

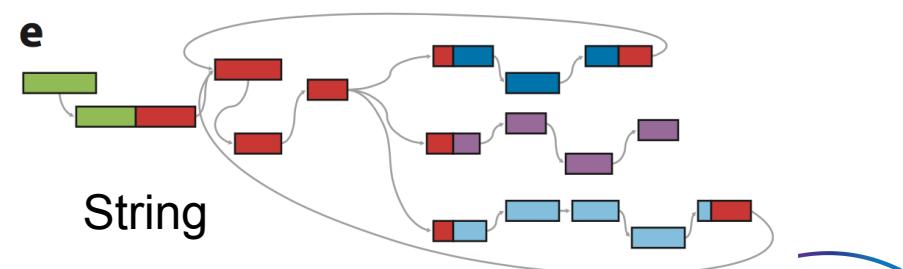
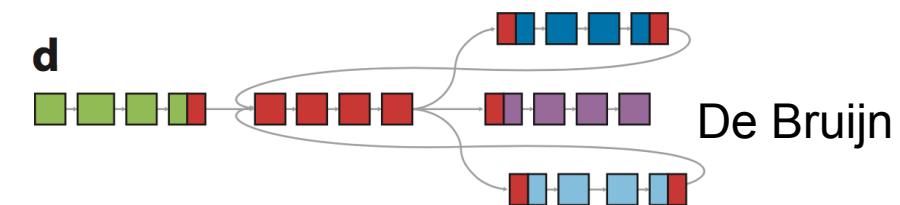
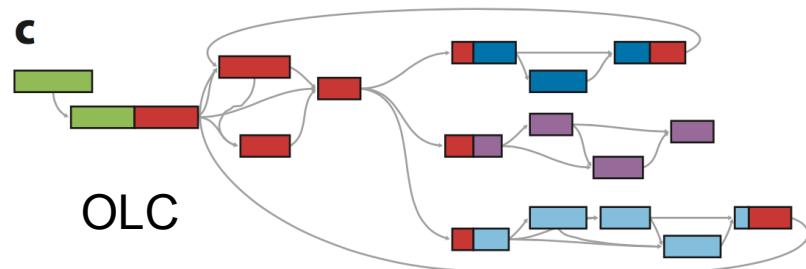
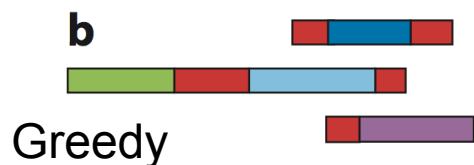
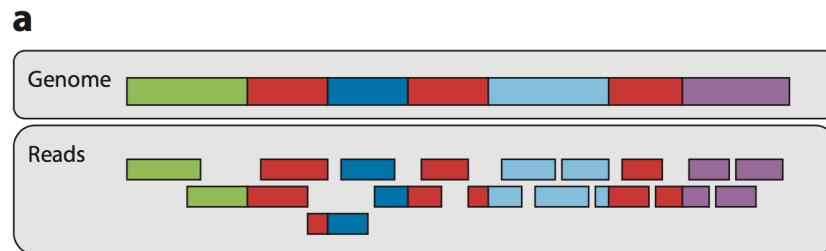


At its core, *de novo* assembly is about looking at how reads (or parts of reads) overlap with each other to determine their relative order and orientation with respect to the original genome.



De novo assembly approaches

- Greedy approaches
- Graph-based approaches
 - Overlap-Layout-Consensus (OLC)
 - De Bruijn Graphs



Source: Simpson, Jared T., and Mihai Pop. Annual review of genomics and human genetics 16 (2015): 153-172.

Greedy approaches

- Procedure:
 1. Sequences (reads/contigs)
 2. All-vs-all pairwise comparison
 3. Find pair with best match/overlap
 4. Collapse/assemble pair into a single contig
 5. Repeat
- Very simple
- Makes the most locally optimal (“greedy”) choice at each step
- Used to great success in early sequencing projects like the Human Genome Project

Greedy approaches

Exercise: Assemble these four contigs using a greedy approach.

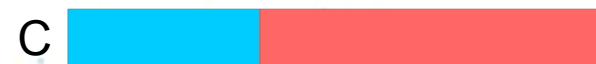
A



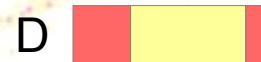
B



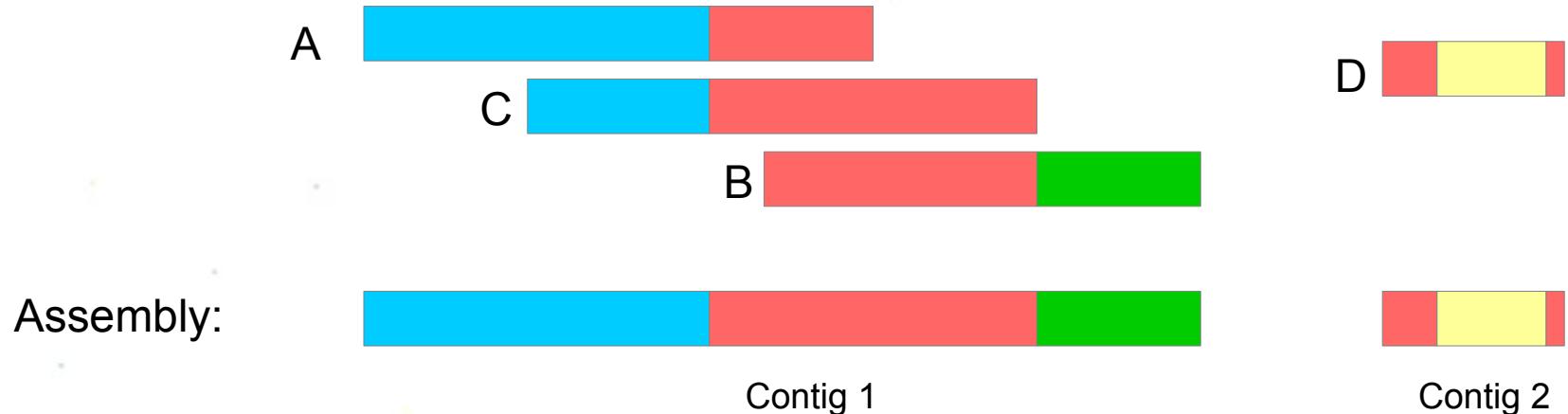
C



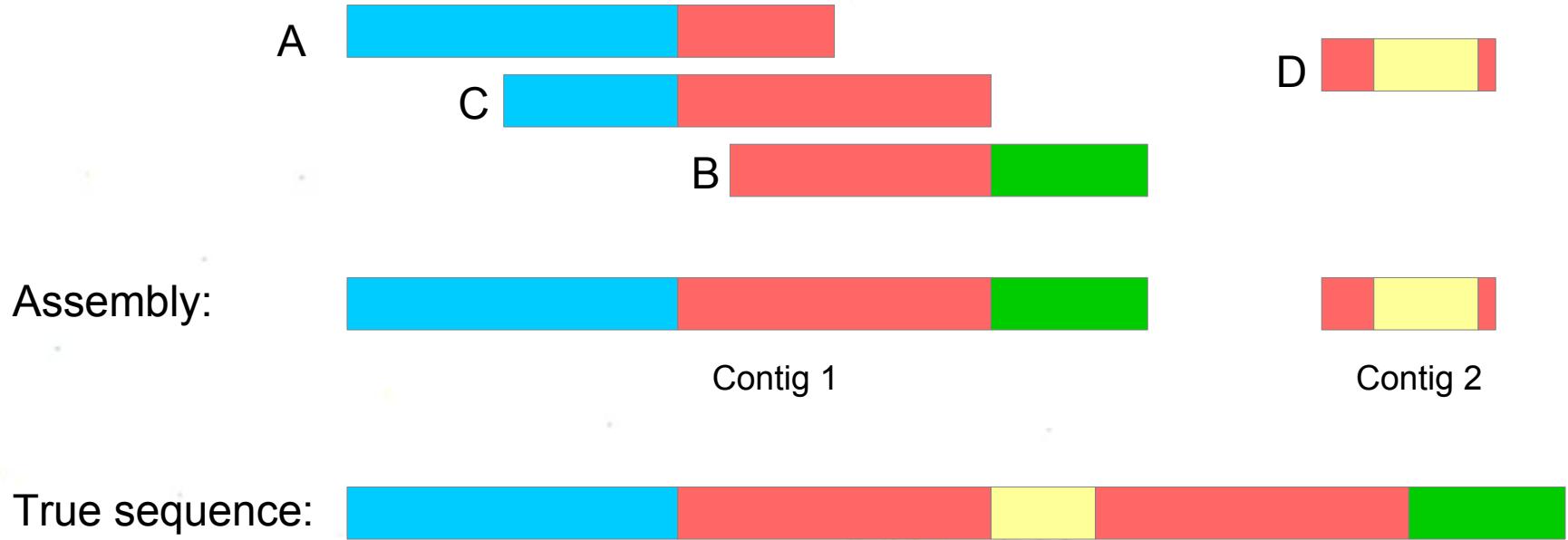
D



Greedy approaches

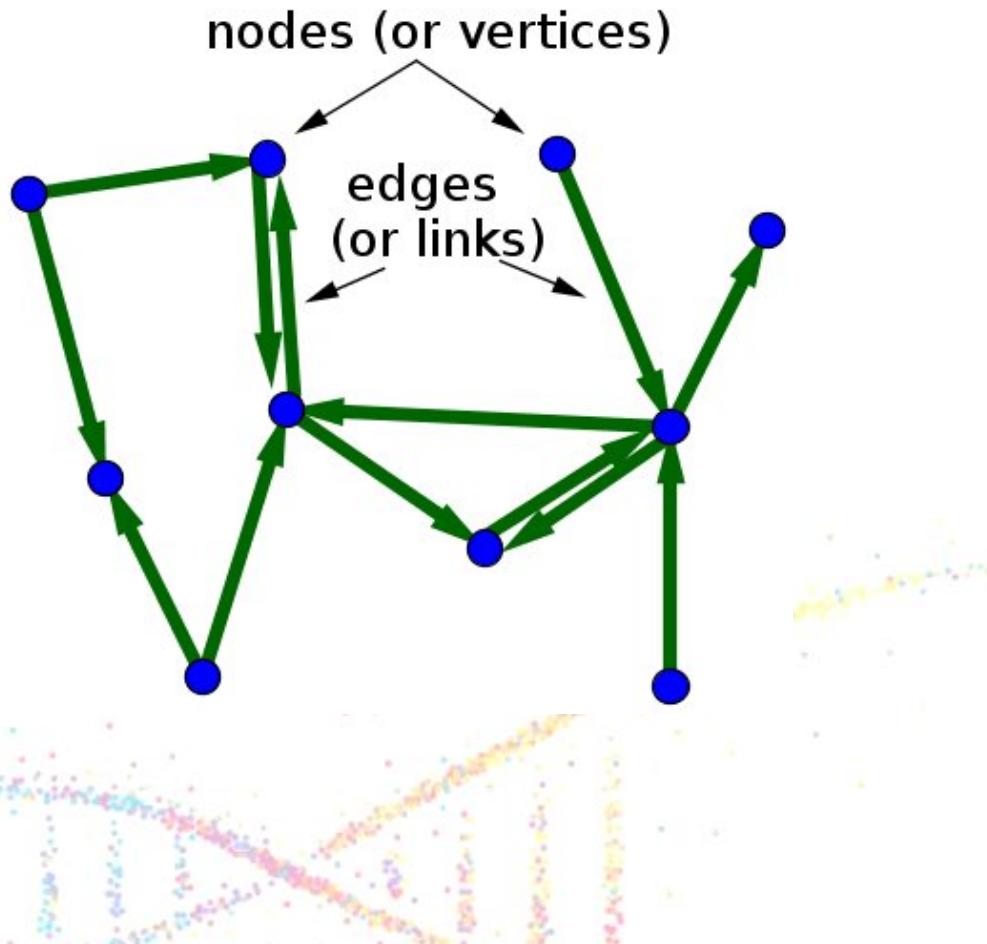


Greedy approaches



- Repeated genomic regions aren't handled well
- Paired-end information can't be easily used
- NGS produces too many reads to compare all vs all
- Now generally replaced by more complex graph-based algorithms

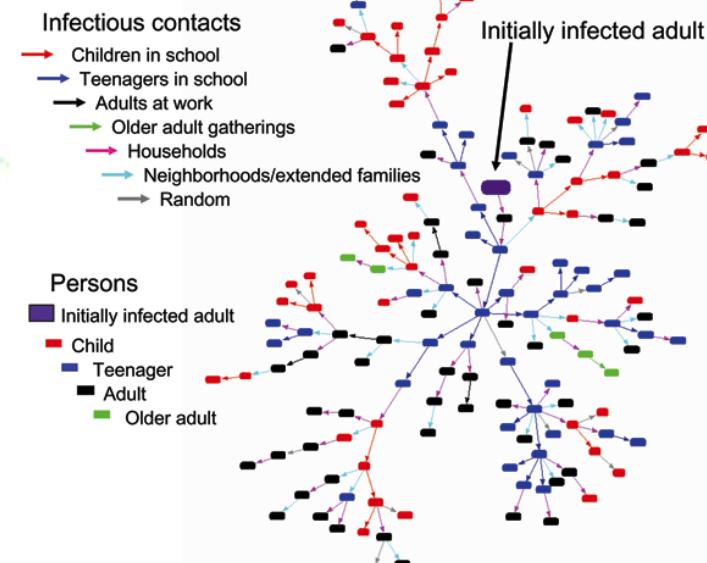
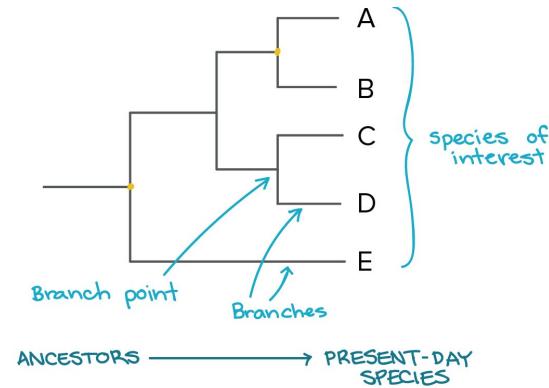
Graphs



http://mathinsight.org/definition/directed_graph

<https://www.khanacademy.org/science/biology/her/tree-of-life/a/phylogenetic-trees>

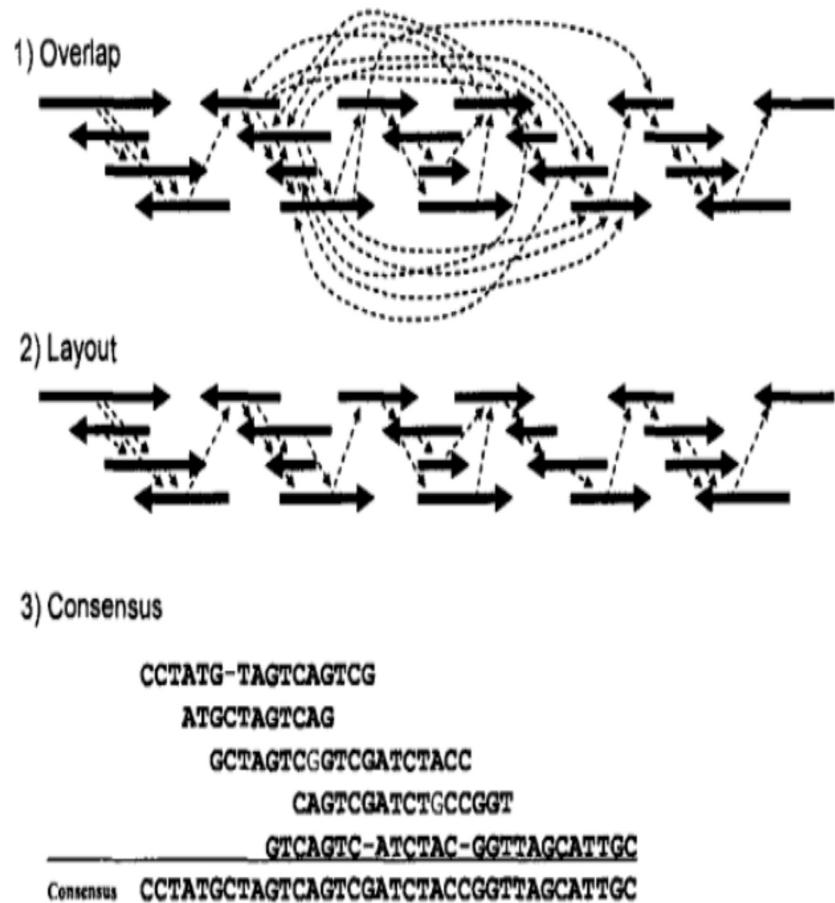
<https://wwwnc.cdc.gov/eid/article/12/11/06-0255-f4>



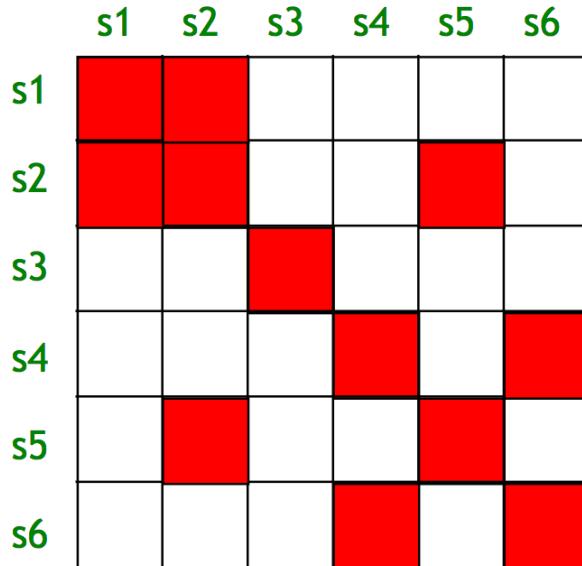
Overlap-Layout-Consensus

Each sequence is a vertex in the graph. Vertices are linked with an edge if they overlap.

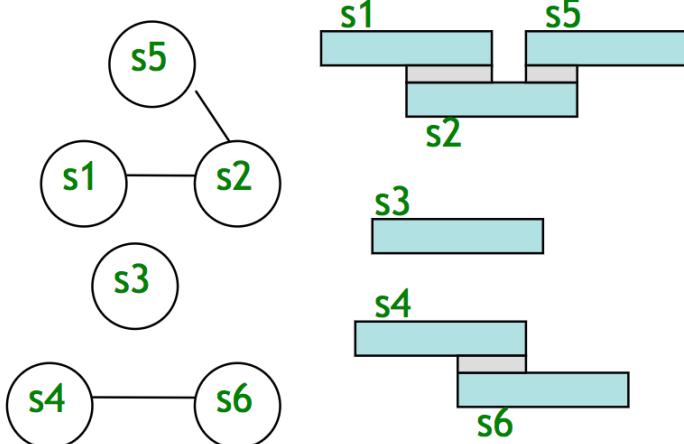
1. Find the **overlaps** by aligning all the reads with each other.
2. **Layout** the reads based on which aligns to which.
3. Join all overlapping read sequences to get the **consensus**.



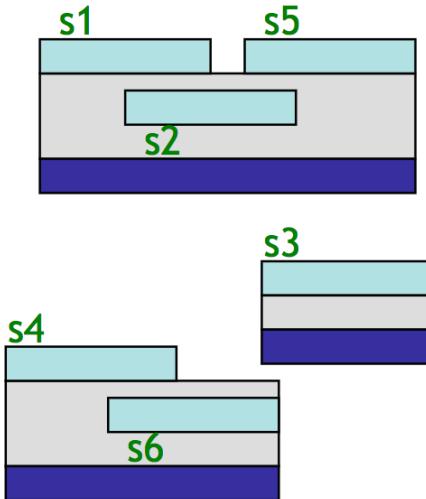
overlap



layout



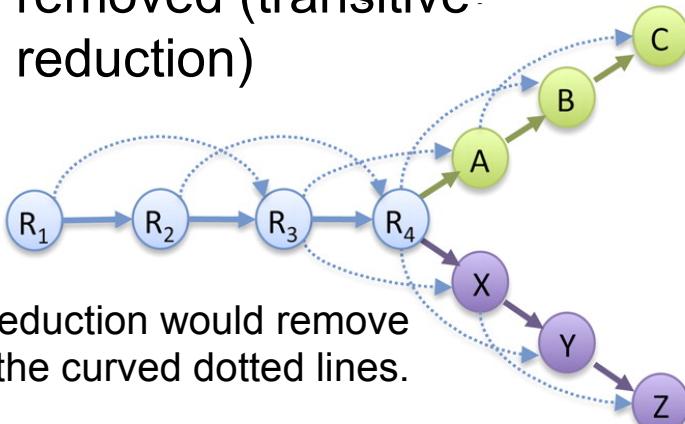
consensus



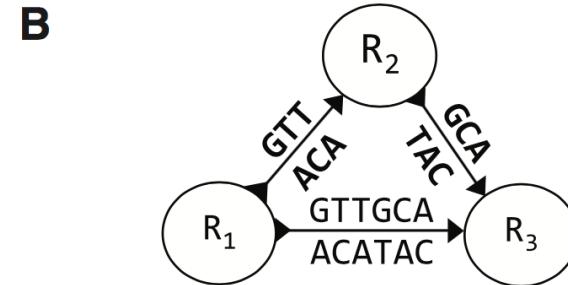
- Like greedy algorithms, OLC requires all-vs-all pairwise comparisons, which is computationally expensive.
- Finding the globally optimum solution during layout is generally impossible. Progressive pairwise alignment needed to estimate.
- Resolves repeats fairly well, using mate-pair constraints.
- Used to great success with earlier genome sequencing projects.
- Generally unsuited to modern short read NGS technologies.

String graphs

- String graphs are made from overlap graphs, with two modifications:
 - Contained reads – reads that are substrings of some other read – are removed
 - Transitive edges are removed (transitive reduction)

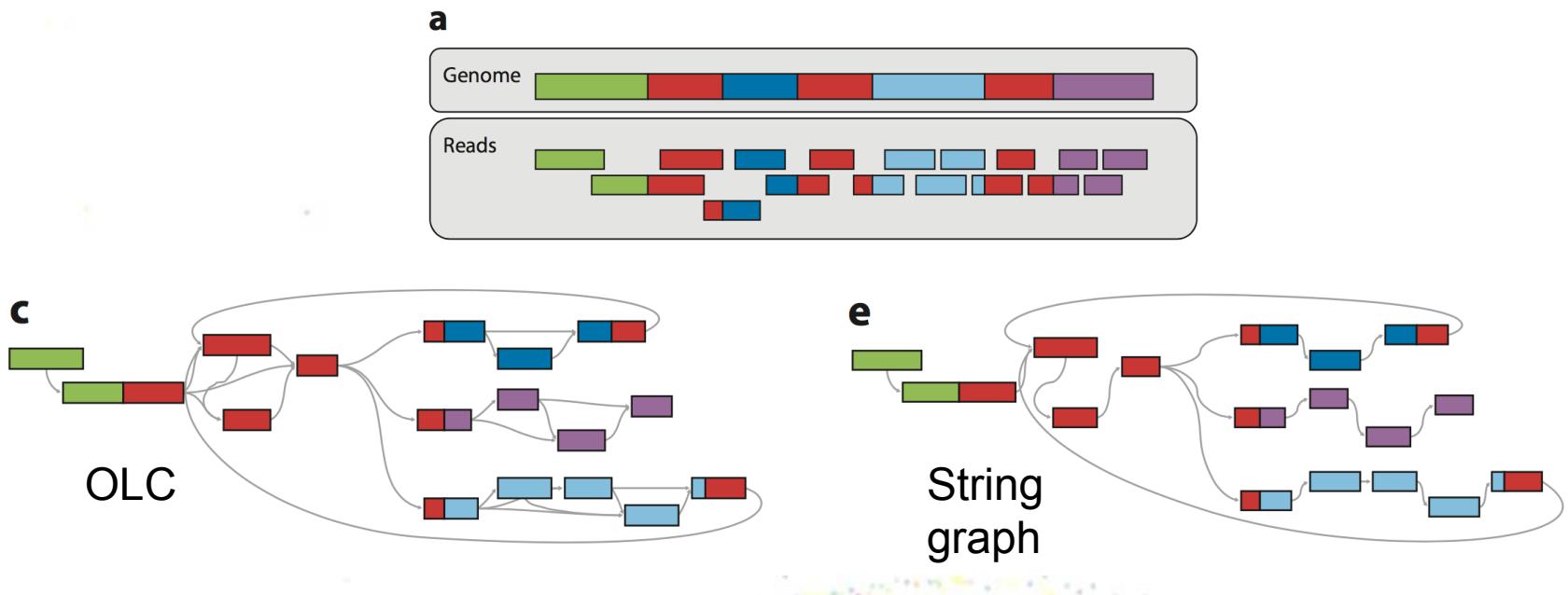


A R_1 ACATACGATACA
 R_2 TACGATACAGTT
 R_3 GATACAGTTGCA



The edge $R_1 \rightarrow R_3$ is redundant with the path $R_1 \rightarrow R_2 \rightarrow R_3$. The edge $R_1 \rightarrow R_3$ is therefore *transitive*, and can be removed.

OLC vs String graphs



String graphs represent the structure of the genome in a more concise way, addressing the computational complexity of the OLC approach. This has lead to new, more memory-efficient assemblers compatible with short read NGS technology.

De Bruijn Graphs

A de Bruijn graph models the relationship between exact substrings of length k , called k -mers, extracted from input reads. Vertices in the graph represent k -mers, and edges indicate that two k -mers overlap by exactly $k-1$ bases.

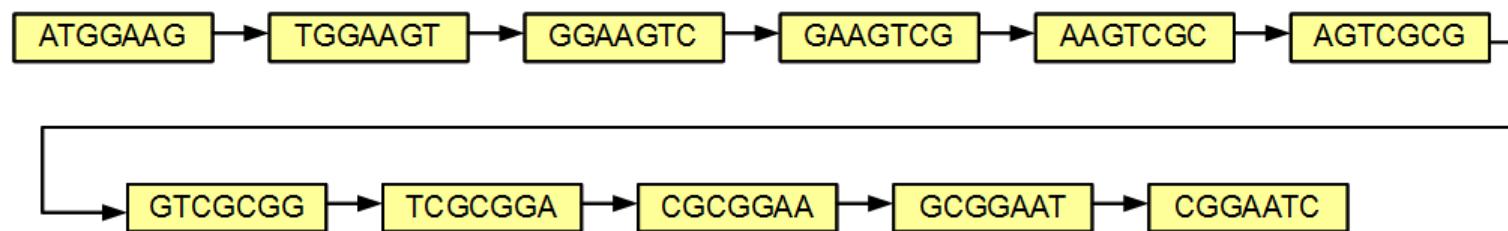
sequence

ATGGAAGTCGCGGAATC

7mers

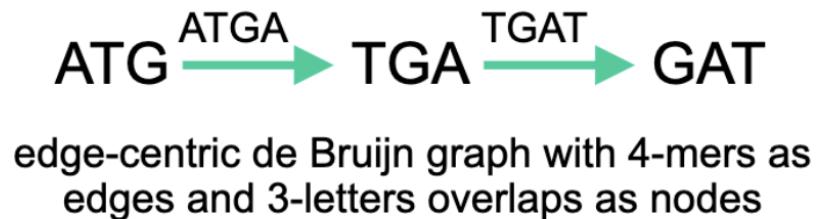
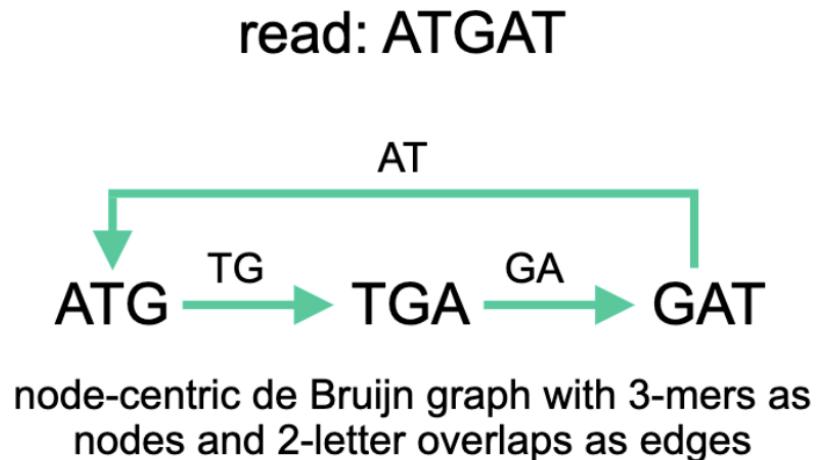
ATGGAAG
TGGAAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC

de Bruijn graph



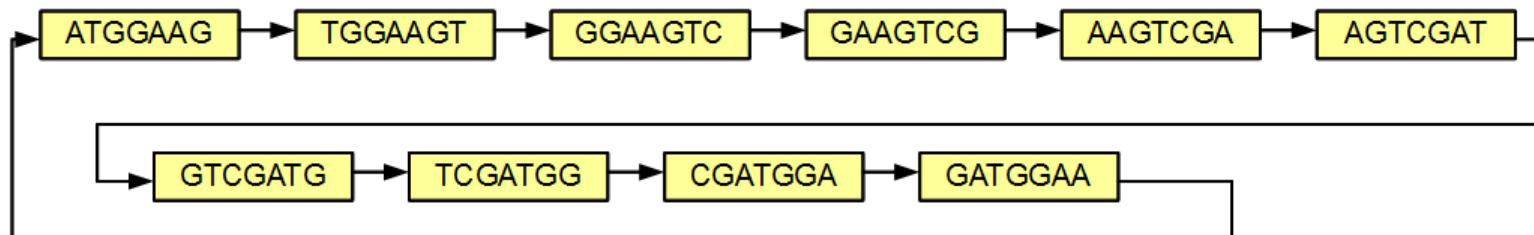
De Bruijn Graphs

- Node-centric
 - The nodes are all the unique k -length substrings in the read set, and an edge is created if there is a $(k-1)$ -length overlap.
- Edge-centric
 - The nodes are the $(k-1)$ -length substrings, and an edge connects two nodes if their $(k-1)$ -mers are consecutive in some read.

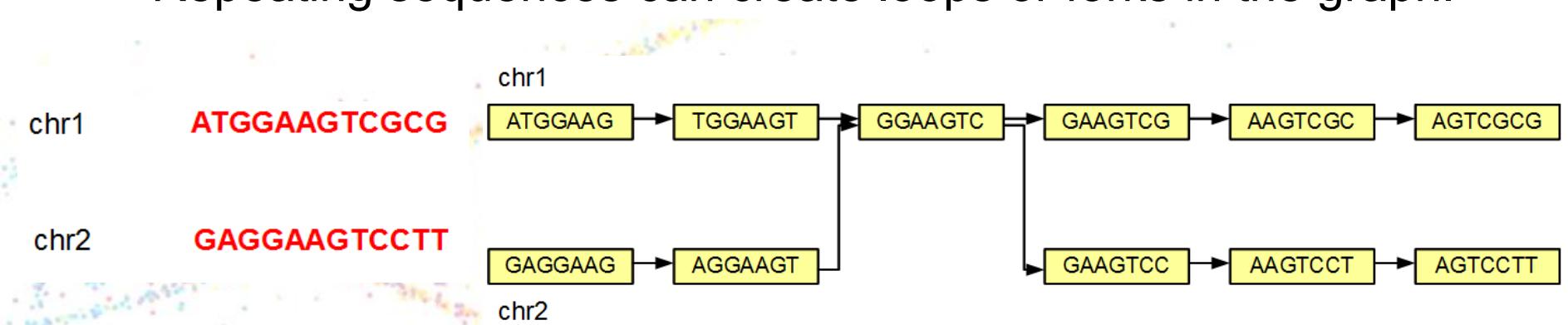


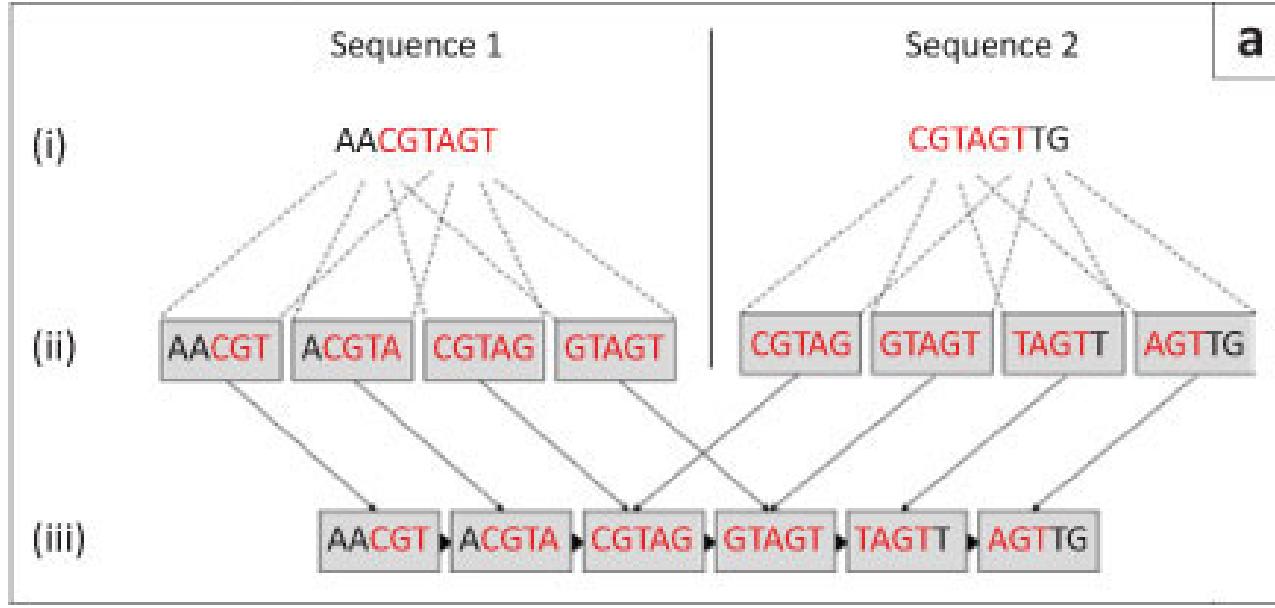
ATGGAAGTCGATGGAAAG

ATGGAAG
TGGAAGT
GGAAGTC
GAAGTCG
AAGTCGA
AGTCGAT
GTCGATG
TCGAIGG
CGATGGA
GATGGAA
ATGGAAG



Repeating sequences can create loops or forks in the graph.

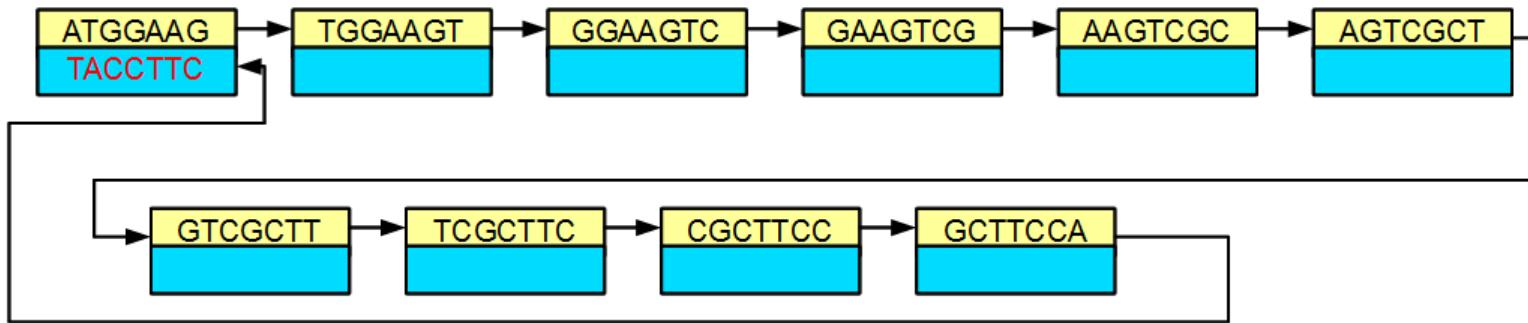




By breaking down reads into k -mers, de Bruijn graphs allow overlapping reads to be correctly linked and ordered, without having to compute an overlap score.

ATGGAAAGTCGCTTCCAT

ATGGAAAG
TGGAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCT
GTCGCTT
TCGCTTC
CGCTTCC
GCTTCCA
CTTCCAT



- k -mers can be represented in the de Bruijn graph in double-stranded fashion. An edge can indicate that a k -mer overlaps with the reverse-complement of another.
- k -mers of even length can be reverse-complements of themselves (e.g. ATAT), creating ambiguity in the graph. For this reason, only k -mers of odd length are typically used.

De Bruijn Graphs

Exercise: Make an edge-centric de Bruijn graph for the word “BANANA”. Use $k=4$.

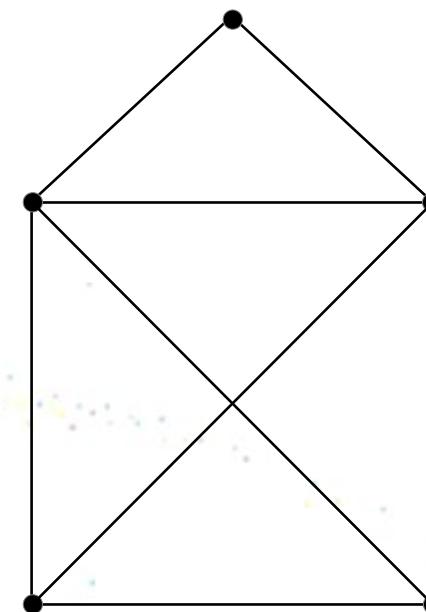
1. Make a list of all the k -mers.
2. Split each k -mer into the left and right $(k-1)$ -mer. Each $(k-1)$ -mer forms a node in your graph.
3. For each k -mer, draw an edge between the left and right $(k-1)$ -mer.

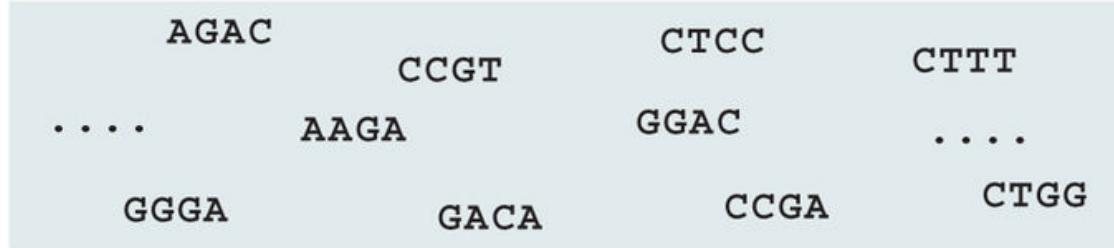


Ok, now you have a de Bruijn graph. How do you use it to assemble the original sequence?

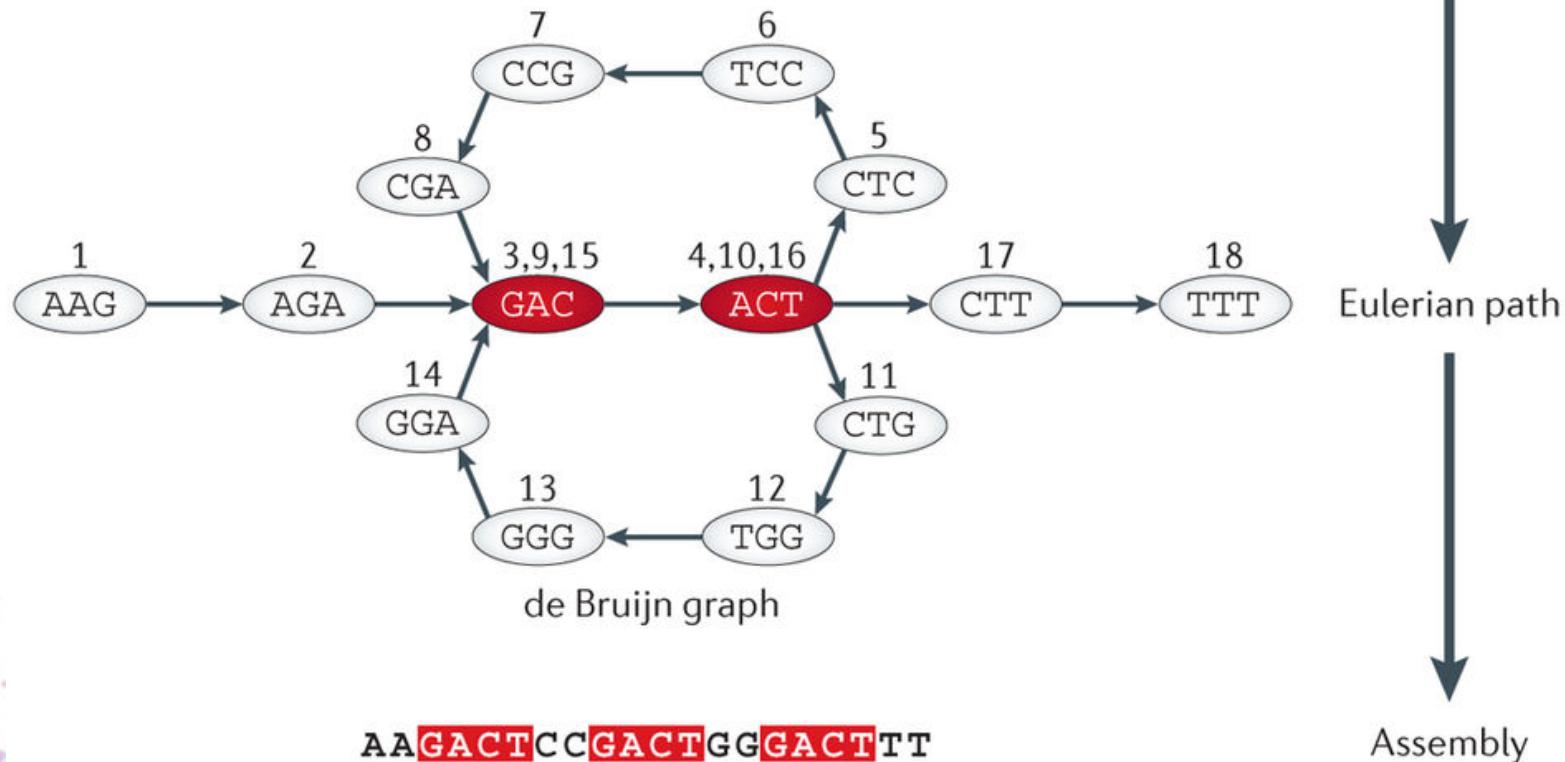
Find the Eulerian path!

Exercise: Draw the following diagram without lifting your pen or tracing over a line more than once.





Reads



The original sequence is reconstructed by following the Eulerian path through the graph, wherein every edge is traversed exactly once.

Unknown target genome

ATGCTATGCGT

reads

ATGCTA
CTATGC
ATGCGT

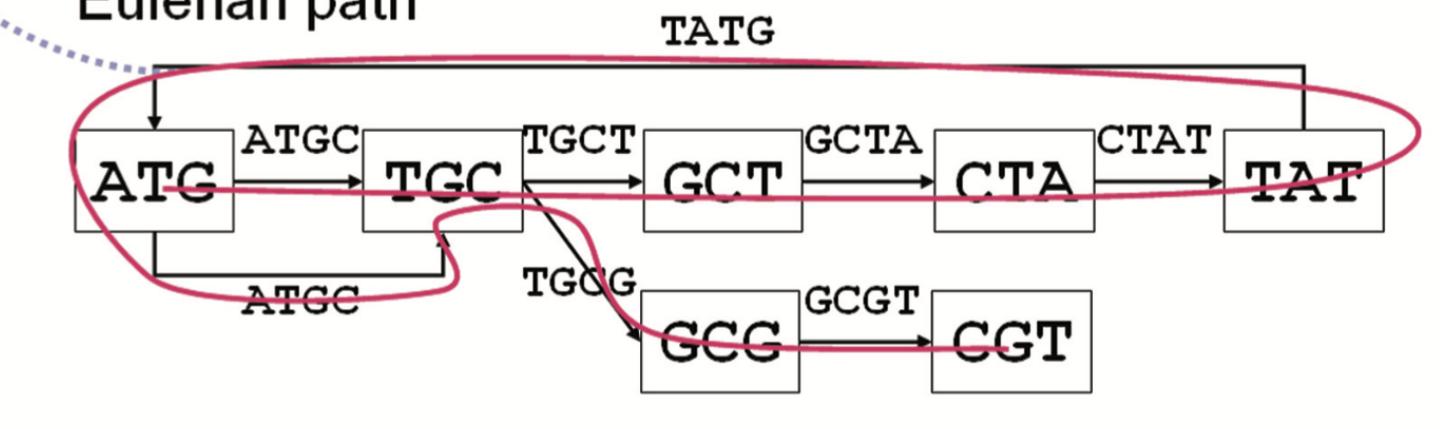


k-mers ($k=4$)
ATGC TGCT GCTA CTAT
TATG ATGC TGCG GCGT

ATGC → ATG, TGC
 k -mer $(k-1)$ -mer

de Bruijn Graph

Eulerian path



Unknown target genome

ATGCTATGCGT

reads

ATGCTA
CTATGC
ATGCGT



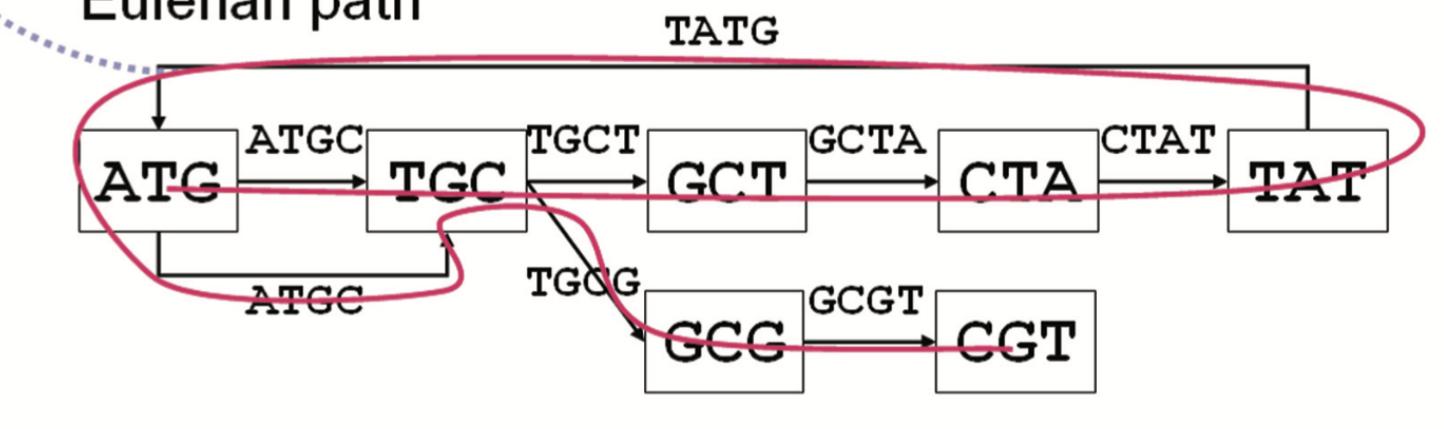
k-mers ($k=4$)

ATGC TGCT GCTA CTAT
TATG ATGC TGCG GCGT

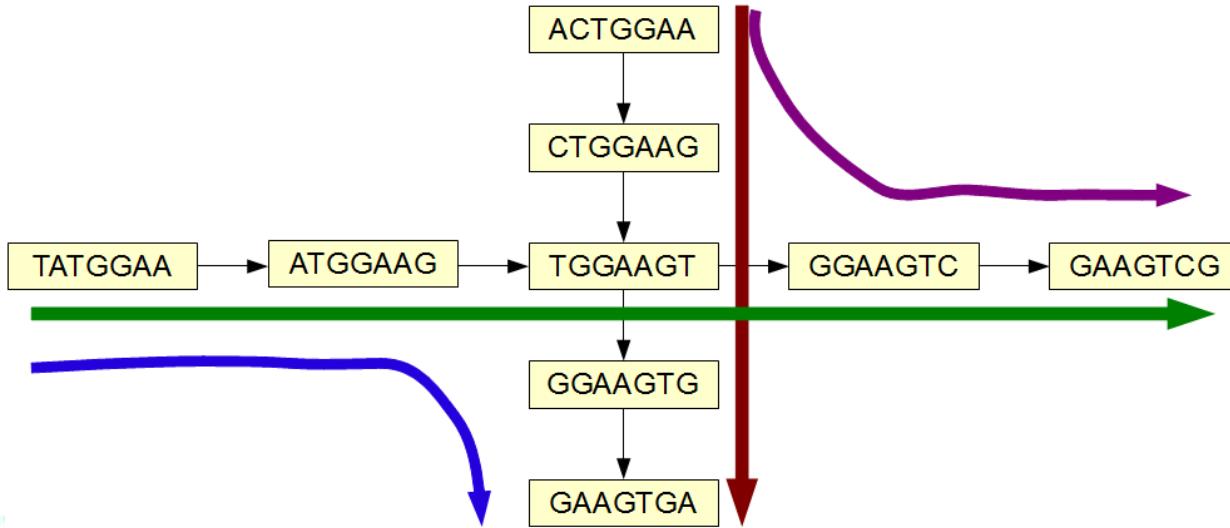
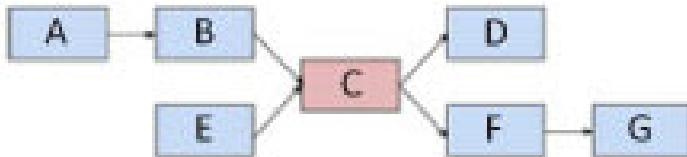
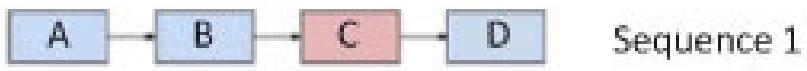
ATGC → ATG, TGC
 k -mer $(k-1)$ -mer

de Bruijn Graph

Eulerian path



Exercise: Draw the de Bruijn graph for the same reads, but with 5-mers ($k=5$).

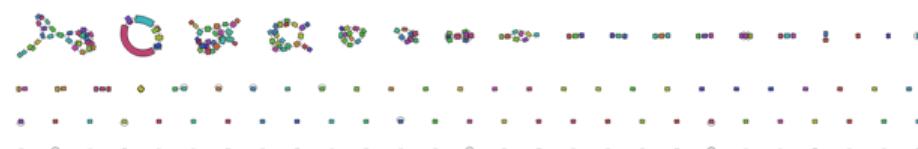
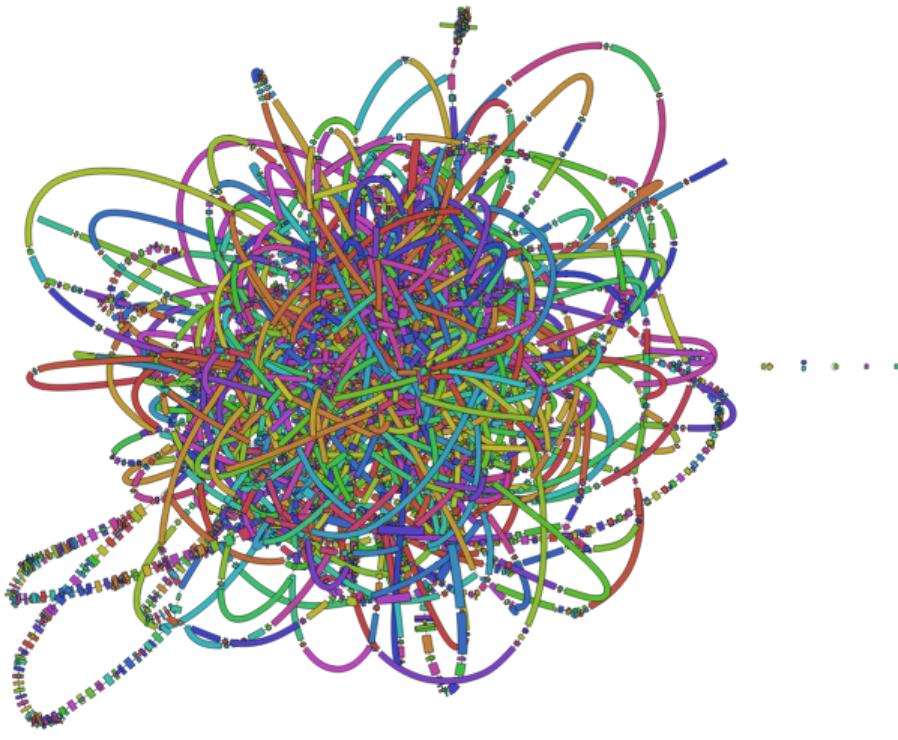


Original sequences – TATGGAAAGTCG, ACTGGAAAGTGA

By breaking down reads into smaller k-mers, de Bruijn graphs lose useful positional information. This can create links between unrelated sequences that cannot be entangled without consulting the original reads.

The choice of k -mer length is important. Here, a longer k -mer would have resolved the ambiguity.

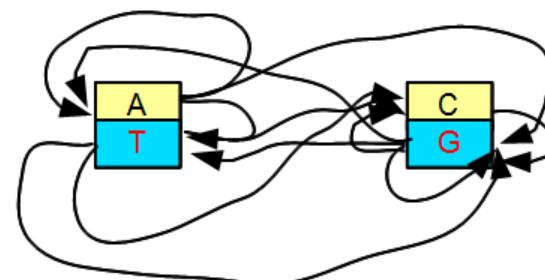
At low values for k , the resulting graph becomes complex and tangled from spurious associations, making final assembly difficult. At sufficiently low values of k , the de Bruijn graph becomes meaningless.

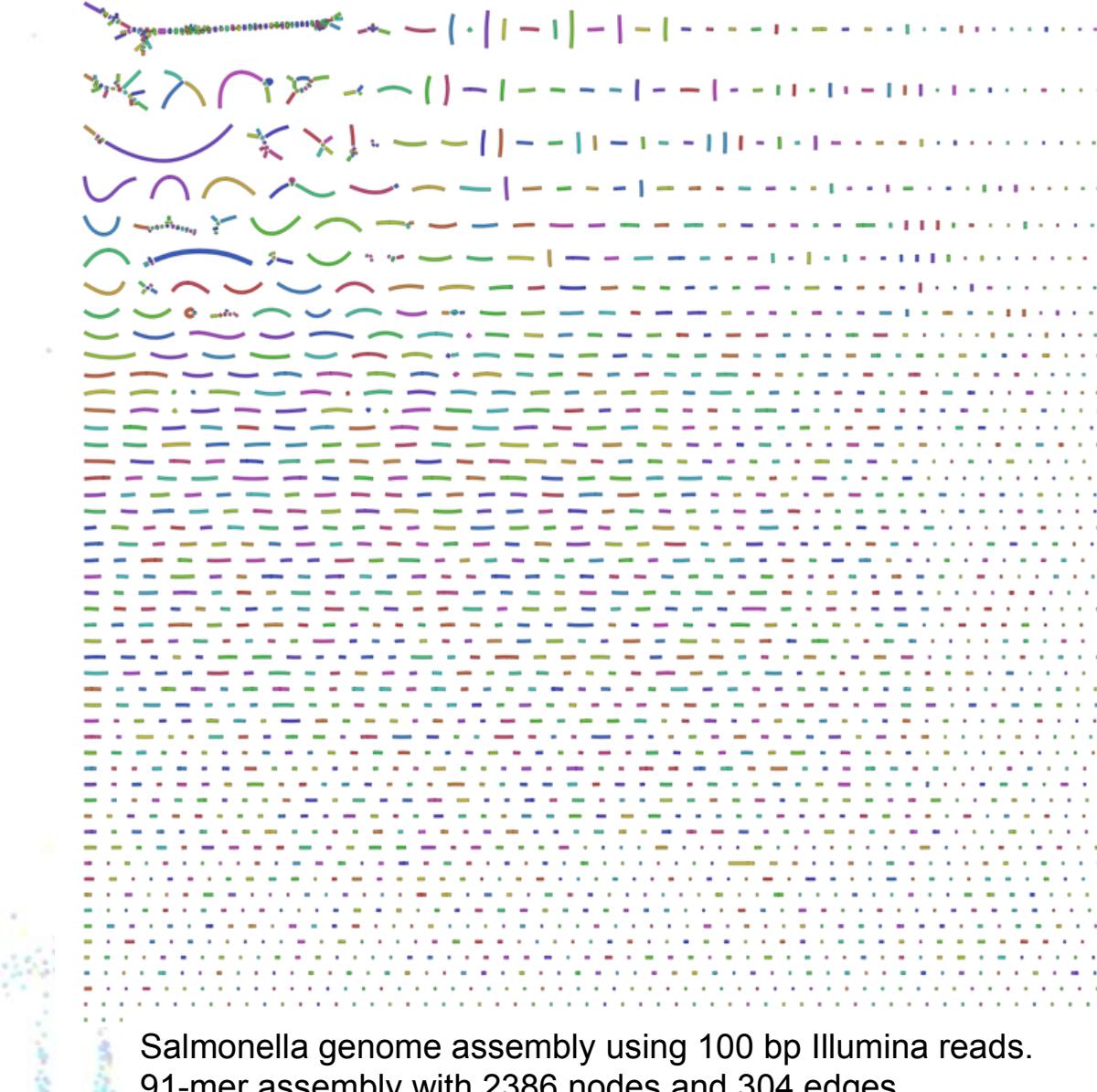


Salmonella genome assembly using 100 bp Illumina reads.
51-mer assembly with 4618 nodes and 6070 edges.

ATGGAAAGTCGCGGAATC

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

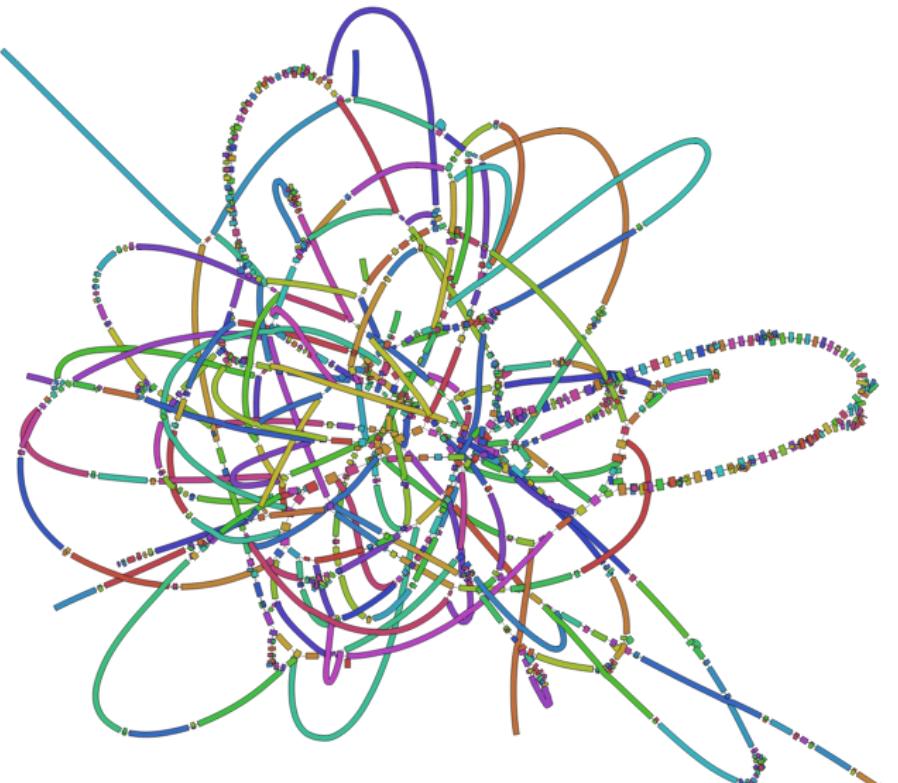




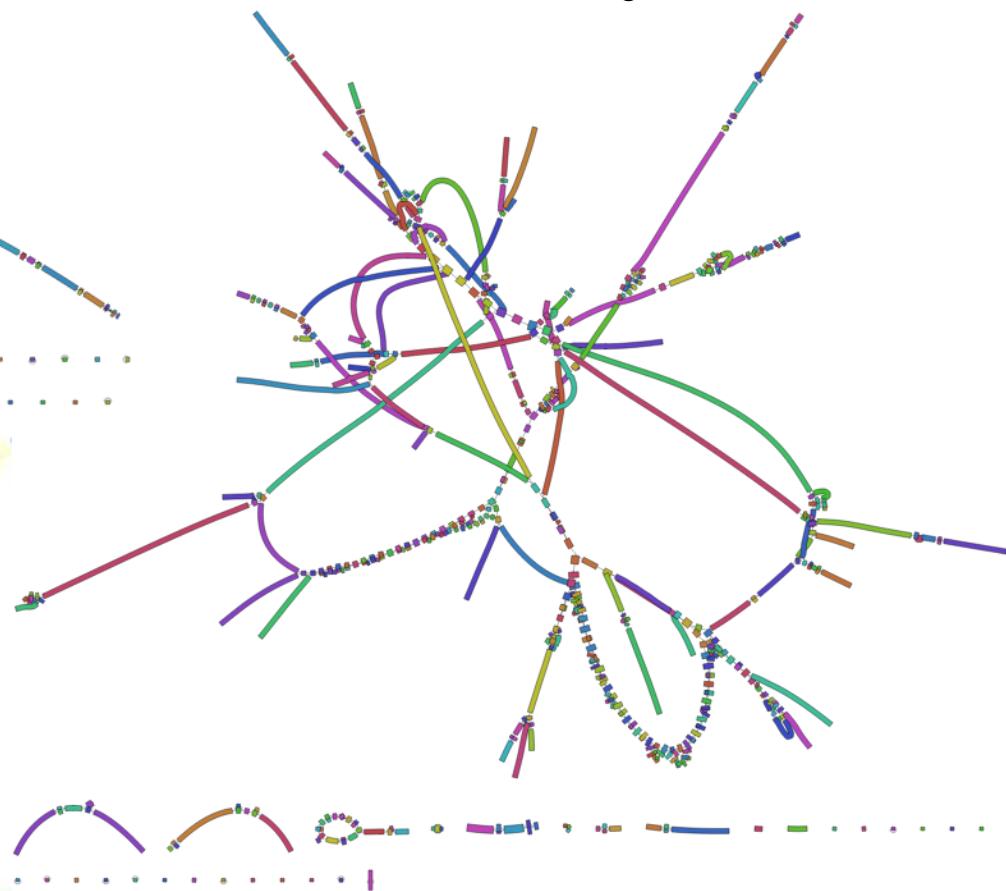
Salmonella genome assembly using 100 bp Illumina reads.
91-mer assembly with 2386 nodes and 304 edges.

Source: <https://github.com/rwick/Bandage/wiki/Effect-of-kmer-size>

At high values of k , you fail to make associations between k -mers from different reads, since the necessary $k-1$ base overlap becomes increasingly unlikely. The resulting de Bruijn graph becomes disjointed and useless.

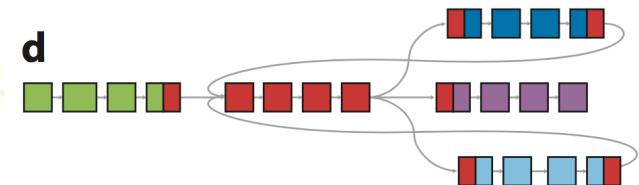
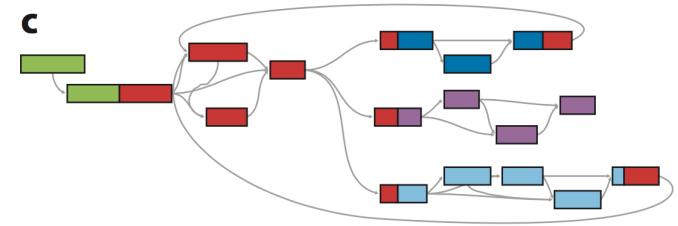


71-mer assembly with 611 nodes and 765 edges.



OLC vs De Bruijn

- Overlap-Layout-Consensus:
 - Appropriate for longer reads
 - Computationally expensive
 - Less sensitive to repeats and errors
 - Key parameter: minimum overlap
- De Bruijn Graphs
 - Appropriate for shorter reads
 - Avoids all-vs-all comparison
 - Repeats get collapsed, simplifying graphs but complicating assembly
 - Computationally more efficient
 - Key parameter: k -mer length



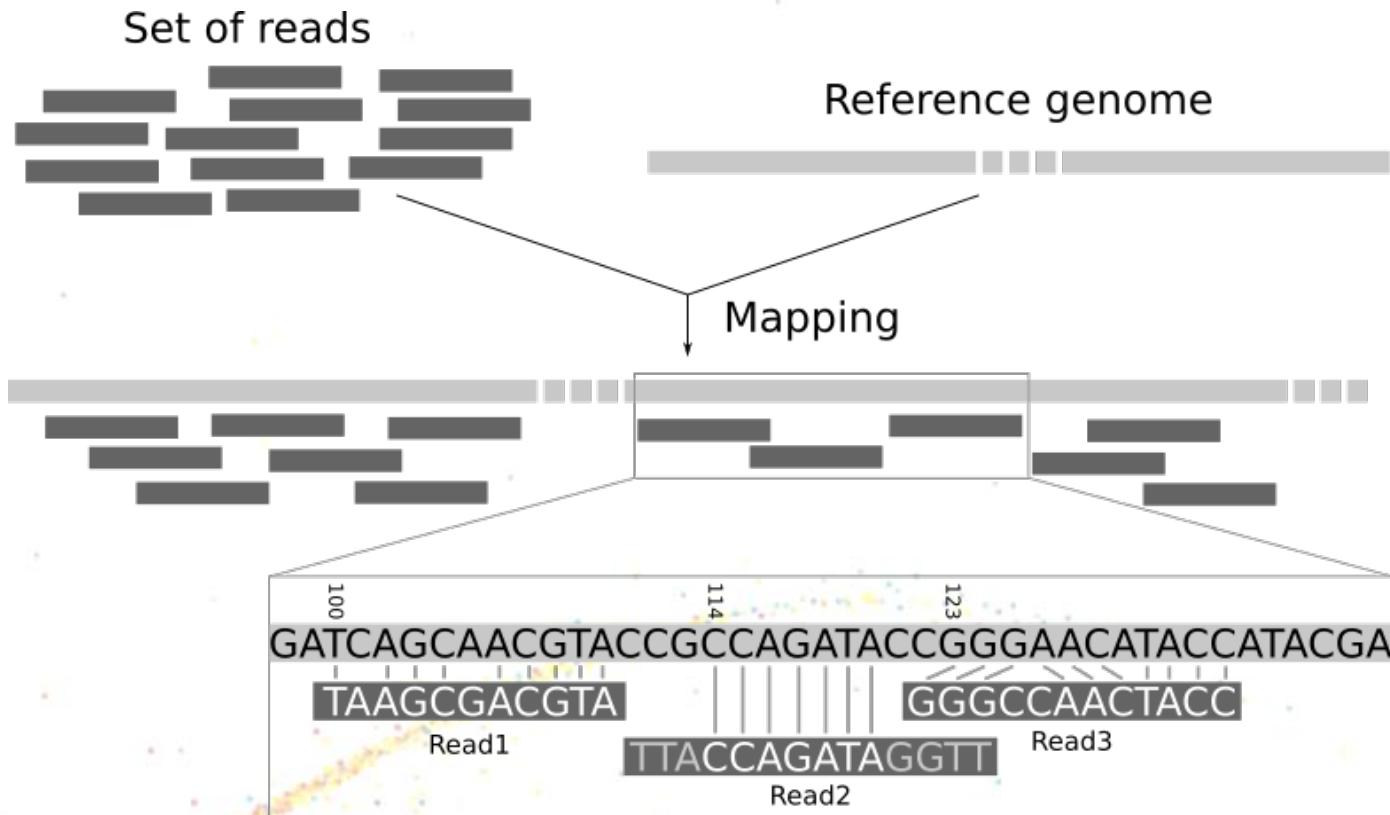


Questions?

Source: Kelly Howe, Lawrence Berkeley Laboratory

Baker, Monya. "De novo genome assembly: what every biologist should know." Nature methods 9.4 (2012): 333-337.

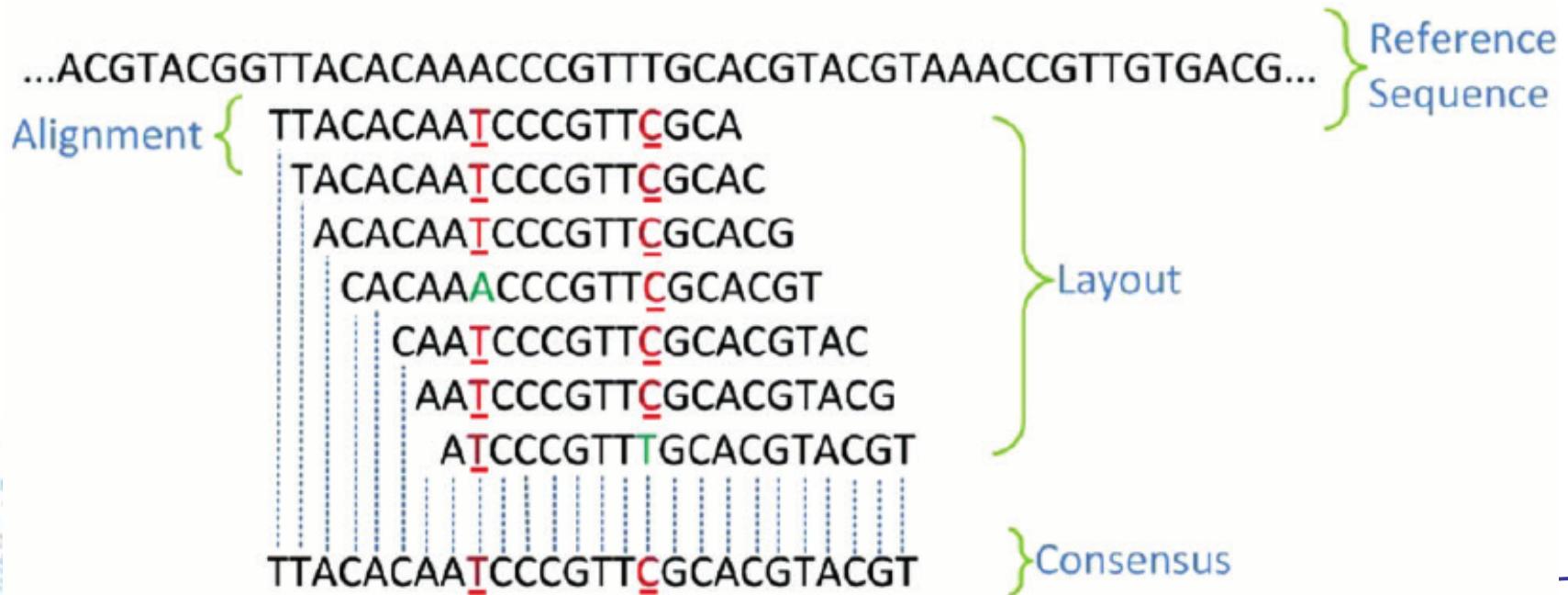
Read mapping



The process of aligning short reads to a longer reference sequence

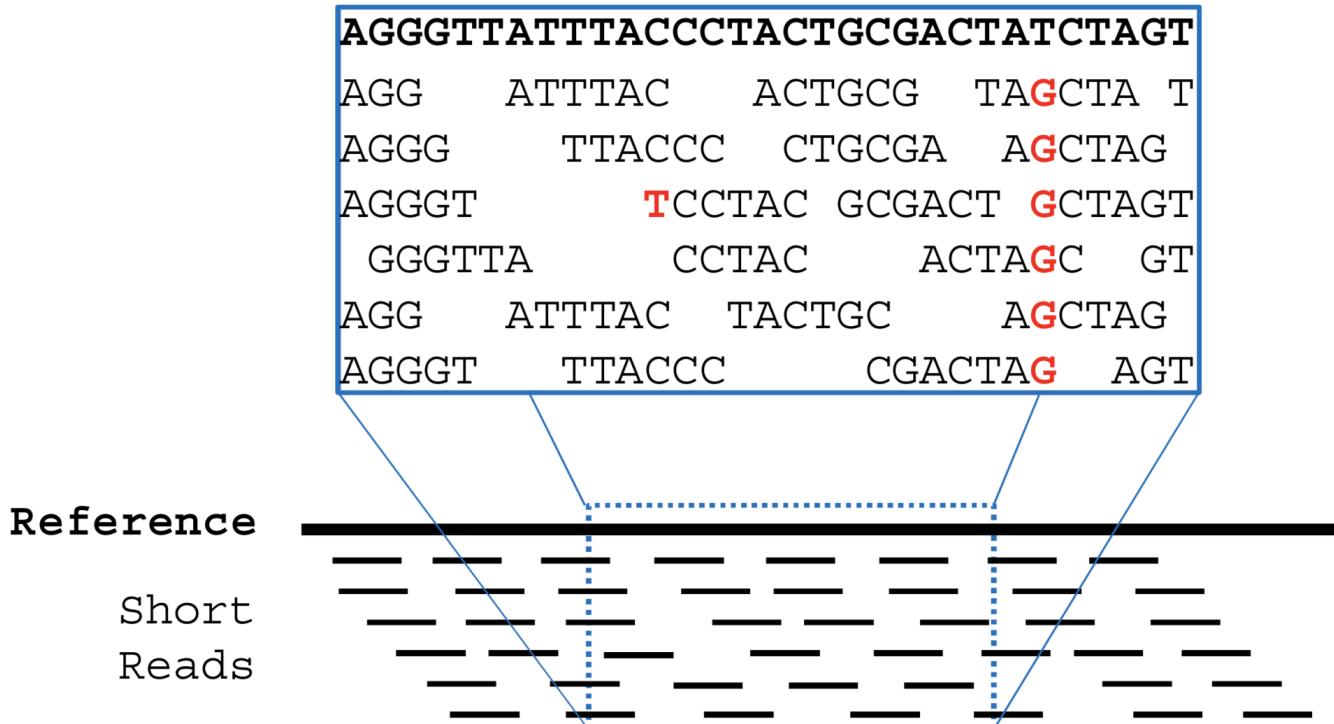
Read mapping

Read mapping is the basis for reference-based genome assembly



Read mapping

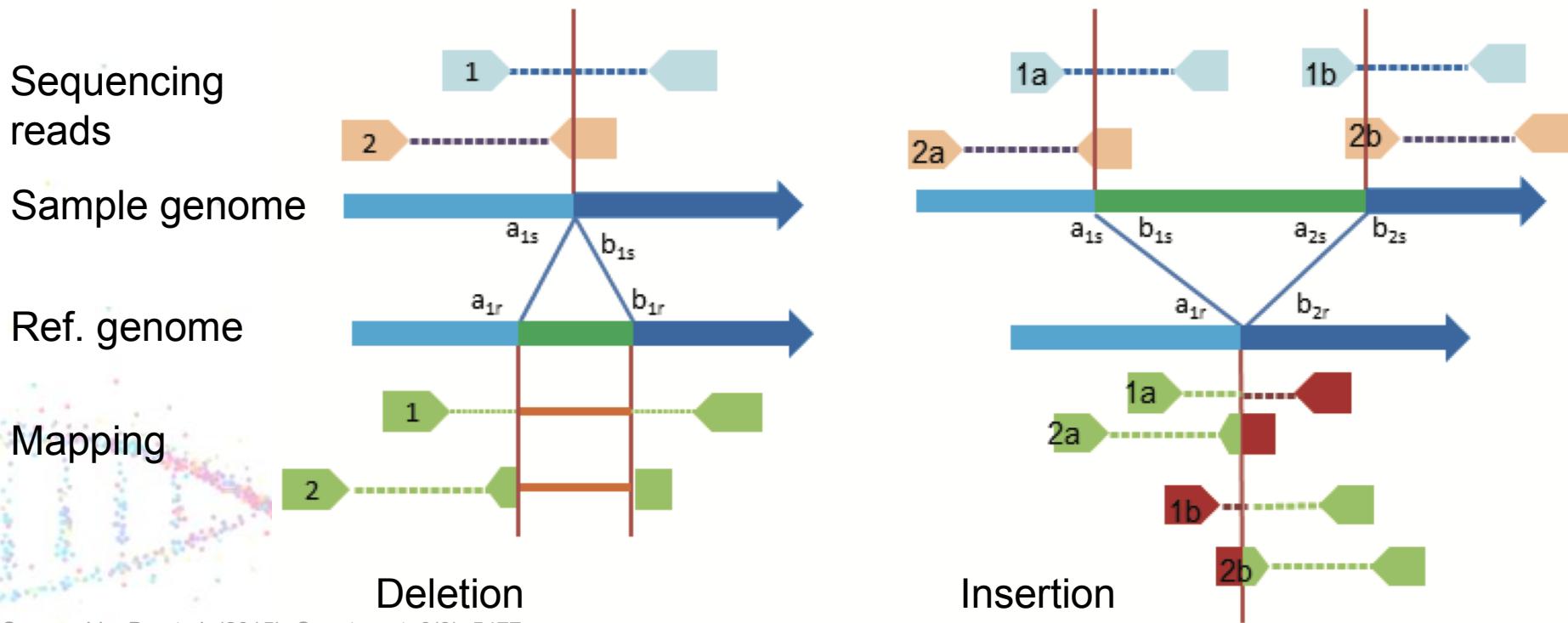
Mapping accuracy is important particularly in identifying structural and single nucleotide variations



Source: Olson, C. B., et al. (2012, April). In 2012 IEEE 20th International Symposium on Field-Programmable Custom Computing Machines (pp. 161-168). IEEE.

Read mapping

Mapping accuracy is important particularly in identifying structural and single nucleotide variations



Read mapping

```
1 agctttcat tctgactgca acgggcaata tgtctctgtg tggattaaaa aaagagtgtc  
61 ttagatgcagg ttctgaactg gttacctgcc gtgagtaaat taaaattta ttgactttagg  
121 tcactaaata cttaaccaa tataggcata gcgcacagac agataaaaaat tacagagtac  
181 acaacatcca tgaaacgcat tagcaccacc attaccacca ccatcaccat taccacaggt  
241 aacggtgcggt gctgacgcgt acaggaaaca cagaaaaaaag cccgcacctg acagtgcggg
```

Shown above are the first 300 bp of the *E. coli* genome.
Which of the following short reads map to reference?

A) catgaaactca

B) ctgatatctgc

C) tgccgtgaaaa

D) caccattacca

Read mapping

```
1 agctttcat tctgactgc a acggcaata tgtctctgtg tggattaaaa aaagagtgtc  
61 ttagatgcagg ttctgaactg gttacctgc gtgagtaaat taaaattta ttgacttagg  
121 tcactaaata cttaaccaa tataggcata gcgcacagac agataaaaaat tacagagtac  
181 acaacatcc tgaaacgcat tagcaccacc attaccaccc ccatcacccat taccacaggt  
241 aacggtgcgg gctgacgcgt acaggaaaca cagaaaaaag cccgcacctg acagtgcggg
```

Shown above are the first 300 bp of the *E. coli* genome.
Which of the following short reads map to reference?

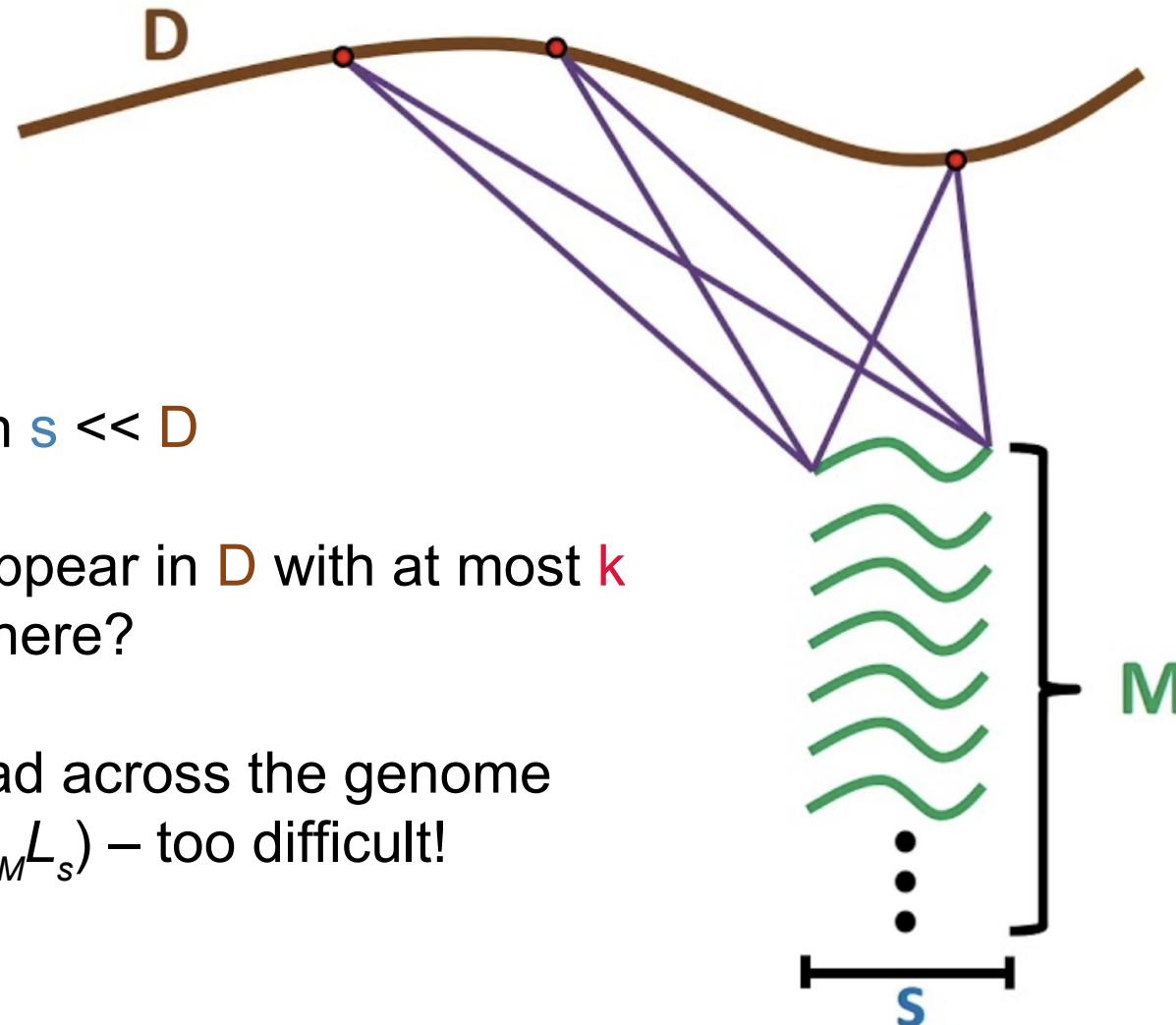
A) catgaaactca

B) ctgatata_{tgc}

C) tgccgtgaaaa

D) caccattacca

Read mapping



Given:

- a genome database **D**
- **M** reads each of length **s** \ll **D**

For each read, does it appear in **D** with at most **k** differences, and if so, where?

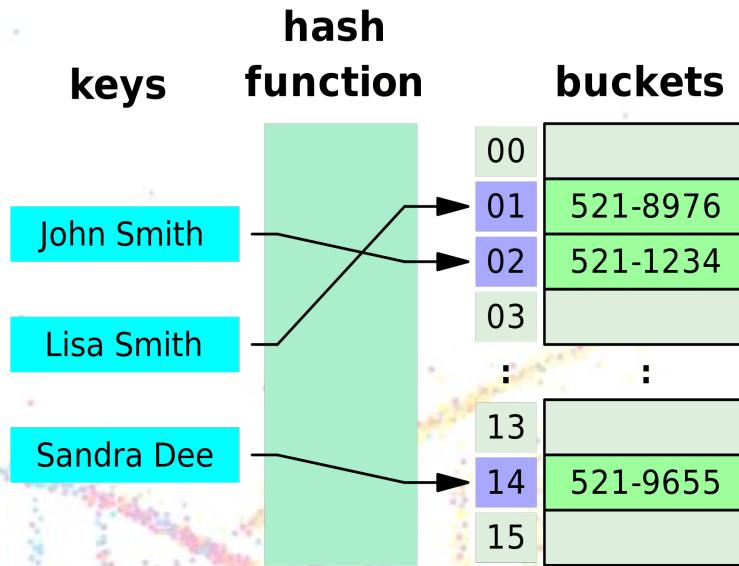
Simply “sliding” each read across the genome has complexity of $O(L_D L_M L_s)$ – too difficult!

Read mapping approaches

- Hash tables
 - MAQ
 - RazerS, RazerS 3
 - RMAP
 - SHRiMP,
SHRiMP2
 - SSAHA, SSAHA2
 - ZOOM
- Burrows-Wheeler Transform
 - BLASR
 - Bowtie, Bowtie 2
 - BWA, BWA-SW,
BWA-MEM
 - SOAP2, SOAP3,
SOAP3-dp

Hash tables

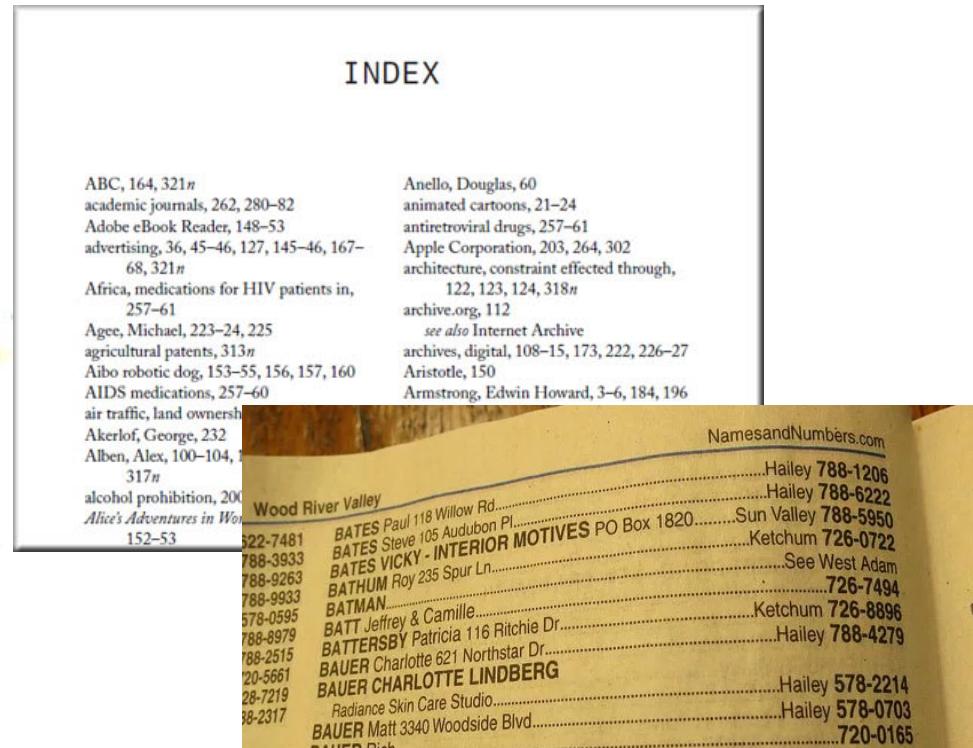
A hash table is a data structure that can match **keys** to **values**. When you input a **key** into the hash function, it quickly returns a unique location or **index** in computer memory where the **value** corresponding to that **key** can be found.



Source: https://en.wikipedia.org/wiki/Hash_table

<https://www.pdfindexgenerator.com/what-is-a-book-index/>

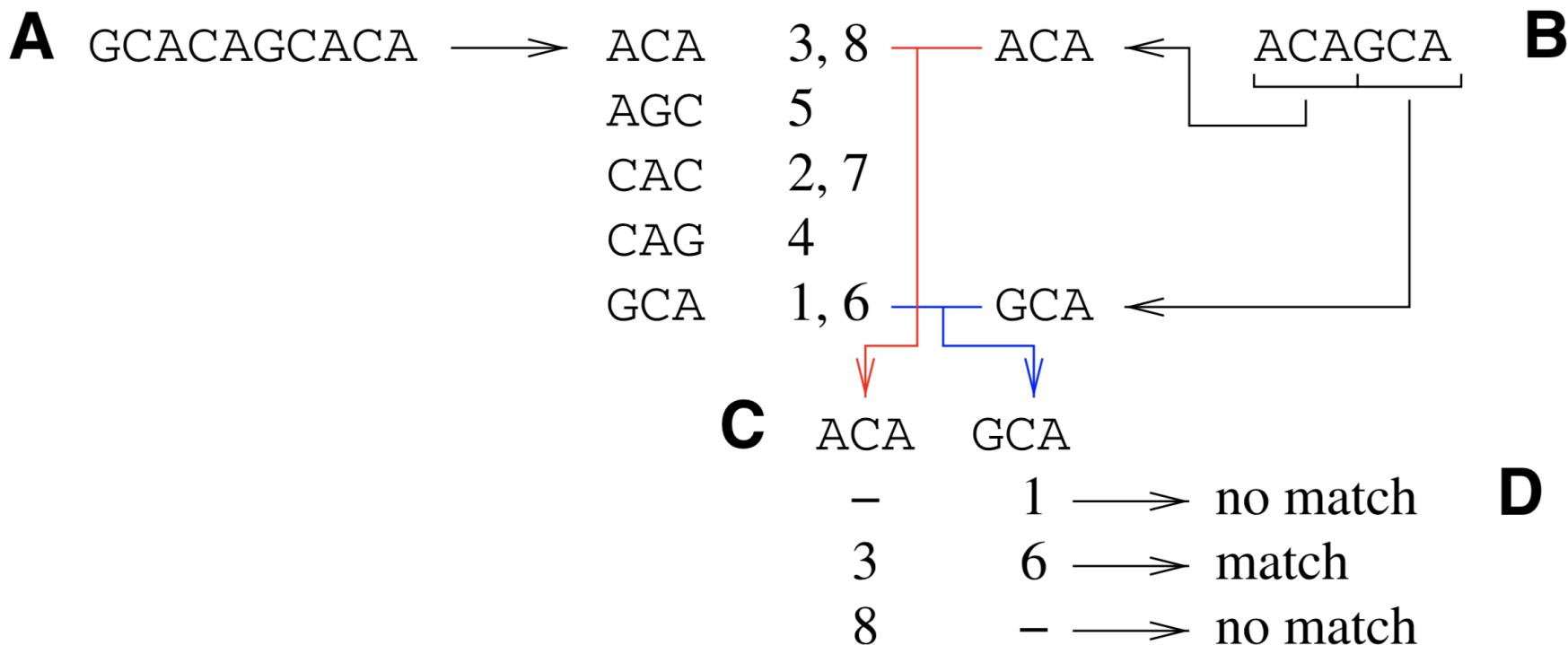
<https://boingboing.net/2017/06/10/how-adam-west-played-a-prank-u.html>



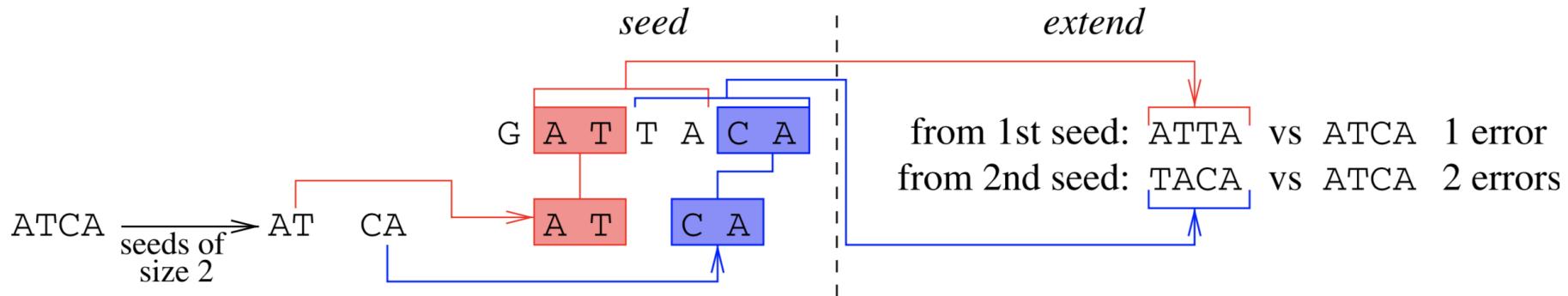
Hash tables

- The “obvious” approach would be to build a hash table mapping all read-length (s) subsequences of the genome to their locations.
 - Requires too much memory and disk space
 - Even harder with variable read lengths
 - Hard to account for mismatches
- Better approach is similar to assembly: break sequences down into k -mers of manageable length.
 - k significantly less than read size, i.e. $k < s$.
 - For example, for $k = 9$, there are at most $4^9 = 262,144$ 9-mers in the genome, regardless of the genome length.

- A. Genome is cut into overlapping 3-mers. Store their respective positions in a hash table.
- B. The read is cut into 3-mers. Look these 3-mers up in the hash table.
- C. Positions for each hit are sorted and compared to those for other hits.
- D. Compatible positions indicate a match between the read and genome.

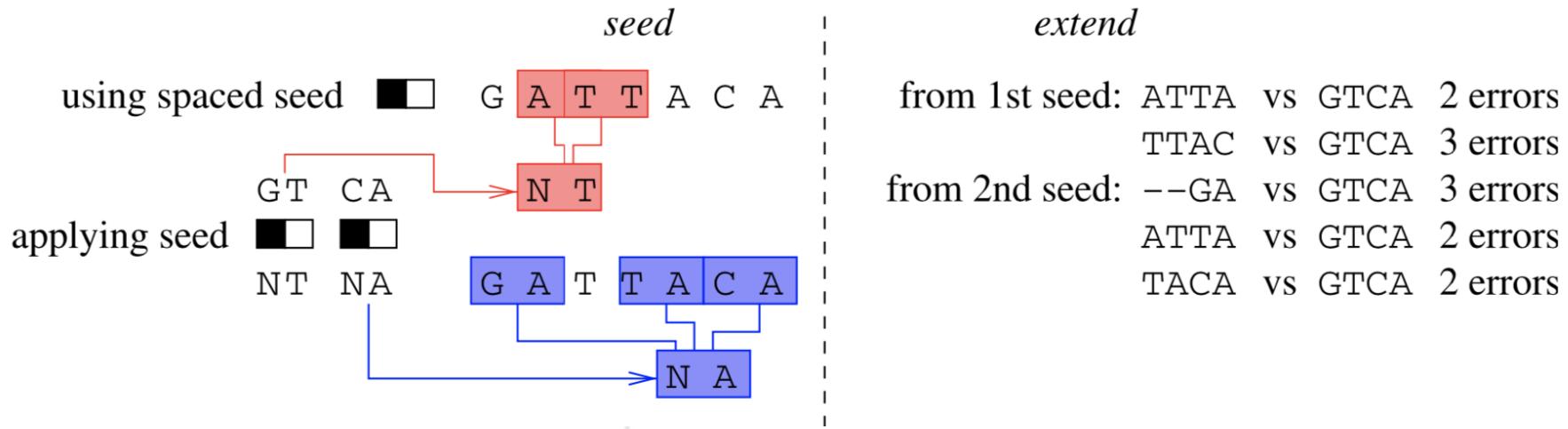


Seed and extend



A read ATCA is sought for in GATTACA, using seeds of size 2, with one error. Each seed maps once (left part). After extension of each seed (right part), it turns out that only one mapping contains just one error.

Seed and extend



As the number of allowed errors increases, the seed must shorten, which can result in too many hits. This is resolved by using “spaced seeds”.

A read GTCA is sought for in GATTACA, using spaced seeds of size 2, with two errors. The first seed maps twice, the second thrice, with one hit in common between them.

Burrows-Wheeler transform

- Invented by Michael Burrows and David Wheeler in 1994 for use in data compression.
- Because of the sorting step, BWT produces runs of repeated characters, making text easy to compress.
- The transformation is easily reversible, with no loss of information.



Source: Schbath, S., et al. (2012). Journal of Computational Biology, 19(6), 796-813.

Burrows-Wheeler transform

genome	suffixes	sorted suffixes	positions	cylinder	suffix array	B-W
GATTACA\$	GATTACA\$	ACA\$	4	ACA\$GATT		T
	ATTACA\$	ATTACA\$	2	ATTACA\$G		G
	TTACA\$	A\$	6	A\$GATTAC		C
	TACA\$	CA\$	5	CA\$GATTA		A
	ACA\$	GATTACA\$	1	GATTACA\$		\$
	CA\$	TACA\$	4	TACA\$GAT		T
	A\$	TTACA\$	3	TTACA\$GA		A
	\$	\$	7	\$GATTACA		A

Step 1: Mark the end of your text (i.e. your genome) with a sigil (\$).

Burrows-Wheeler transform

genome	suffixes	sorted suffixes	positions	cylinder	suffix array	B-W
GATTACA\$	GATTACAS\$	ACA\$	4	ACA\$GATT		T
	ATTACAS\$	ATTACAS\$	2	ATTACA\$G		G
	TTACAS\$	A\$	6	A\$GATTAC		C
	TACAS\$	CA\$	5	CA\$GATTA		A
	ACAS\$	GATTACAS\$	1	GATTACA\$		\$
	CA\$	TACAS\$	4	TACA\$GAT		T
	A\$	TTACAS\$	3	TTACA\$GA		A
	\$	\$	7	\$GATTACA		A

Step 2: List down all the suffixes (right-most substrings) of the text.

Burrows-Wheeler transform

genome	suffixes	sorted suffixes	positions	cylinder	suffix array	B-W
GATTACA\$	GATTACA\$	ACA\$	4	ACA\$GATT		T
	ATTACA\$	ATTACA\$	2	ATTACA\$G		G
	TTACA\$	A\$	6	A\$GATTAC		C
	TACA\$	CA\$	5	CA\$GATTA		A
	ACA\$	GATTACA\$	1	GATTACA\$		\$
	CA\$	TACA\$	4	TACA\$GAT		T
	A\$	TTACA\$	3	TTACA\$GA		A
	\$	\$	7	\$GATTACA		A

Step 3: Sort all the suffixes lexicographically (i.e. alphabetically). By convention, the sigil comes last in the lexicographic order.

Burrows-Wheeler transform

genome	suffixes	sorted suffixes	positions	cylinder suffix array	B-W
GATTACA\$	GATTACA\$	ACA\$	4	ACA\$GATT	T
	ATTACA\$	ATTACA\$	2	ATTACA\$G	G
	TTACA\$	A\$	6	A\$GATTAC	C
	TACA\$	CA\$	5	CA\$GATTA	A
	ACA\$	GATTACA\$	1	GATTACA\$	\$
	CA\$	TACA\$	4	TACA\$GAT	T
	A\$	TTACA\$	3	TTACA\$GA	A
	\$	\$	7	\$GATTACA	A

Step 4: Keep track of the positions for each suffix in the original text.

Burrows-Wheeler transform

genome	suffixes	sorted suffixes	positions	cylinder suffix array	B-W
GATTACA\$	GATTACA\$	ACA\$	4	ACA\$GATT	T
	ATTACA\$	ATTACA\$	2	ATTACA\$G	G
	TTACA\$	A\$	6	A\$GATTAC	C
	TACA\$	CA\$	5	CA\$GATTA	A
	ACA\$	GATTACA\$	1	GATTACA\$	\$
	CA\$	TACA\$	4	TACA\$GAT	T
	A\$	TTACA\$	3	TTACA\$GA	A
	\$	\$	7	\$GATTACA	A

Step 5: "Complete" each suffix, wrapping around from the beginning of the text as if written on a cylinder.

Burrows-Wheeler transform

genome	suffixes	sorted suffixes	positions	cylinder suffix array	B-W
GATTACA\$	GATTACA\$	ACA\$	4	ACA\$GATT	T
	ATTACA\$	ATTACA\$	2	ATTACA\$G	G
	TTACA\$	A\$	6	A\$GATTAC	C
	TACA\$	CA\$	5	CA\$GATTA	A
	ACA\$	GATTACA\$	1	GATTACA\$	\$
	CA\$	TACA\$	4	TACA\$GAT	T
	A\$	TTACA\$	3	TTACA\$GA	A
	\$	\$	7	\$GATTACA	A

Step 6: The last column of the resulting matrix is the Burrows-Wheeler transform (BWT) of the text.

Burrows-Wheeler transform

- The transformation is easily reversible, with no loss of information.
 - Iterative sorting and concatenating of CWT returns the cylinder suffix array

B-W	sorted	concat. 1	sort	concat. 2	sort
T	A	TA	AC	TAC	ACA
G	A	GA	AT	GAT	ATT
C	A	CA	A\$	CA\$	A\$A
A	C	AC	CA	ACA	CA\$
\$	G	\$G	GA	\$GA	GAT
T	T	TT	TA	TTA	TAC
A	T	AT	TT	ATT	TTA
A	\$	A\$	\$A	A\$A	\$GA

- Exercise: which row corresponds to the original string?

Burrows-Wheeler Transform

Exercise: Find the BWT for the word “BANANA”.

1. Mark the end of your text with a sigil (\$).
2. List down all the suffixes of the text.
3. Sort all the suffixes lexicographically.
4. Keep track of the positions for each suffix in the original text (optional).
5. “Complete” each suffix.
6. Extract the last column of the resulting matrix.



Finding a word with BWT

B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted
T	A	T	A	T	A	T	A	T	A
G	A	G	A	G	A	G	A	G	A
C	A	C	A	C	A	C	A	C	A
A	C	A	C	A	C	A	C	A	C
\$	G	\$	G	\$	G	\$	G	\$	G
T	T	T	T	T	T	T	T	T	T
A	T	A	T	A	T	A	T	A	T
A	\$	A	\$	A	\$	A	\$	A	\$

Goal: Find the read “GAT” in the genome “GATTACA”.

Step 1: Starting with the BWT, sort lexicographically, reconstructing the first column of the BWT matrix.

Finding a word with BWT

B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted
T	A	T	A	T	A	T	A	T	A
G	A	G	A	G	A	G	A	G	A
C	A	C	A	C	A	C	A	C	A
A	C	A	C	A	C	A	C	A	C
\$	G	\$	G	\$	G	\$	G	\$	G
T	T	T	T	T	T	T	T	T	T
A	T	A	T	A	T	A	T	A	T
A	\$	A	\$	A	\$	A	\$	A	\$

Goal: Find the read “GAT” in the genome “GATTACA”.

Step 2: We proceed right-to-left in the query string, starting with “T” in “GAT”. Look for all instances of this letter in the sorted BWT.

Finding a word with BWT

B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted
T	A	T	A	T	A	T	A	T	A
G	A	G	A	G	A	G	A	G	A
C	A	C	A	C	A	C	A	C	A
A	C	A	C	A	C	A	C	A	C
\$	G	\$	G	\$	G	\$	G	\$	G
T	T	T	A	T	T	T	T	T	T
A	\$	A	\$	T	A	T	T	A	T
A				A	\$	A	\$	A	\$

Goal: Find the read “GAT” in the genome “GATTACA”.

Step 3: Remember that the BWT column represents the letters immediately preceding the sorted BWT in the text, because the text “wraps” around the cylinder. So let’s look at the corresponding letters in the BWT.

Finding a word with BWT

B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted	B-W	sorted
T	A	T	A	T	A	T	A	T	A
G	A	G	A	G	A	G	A	G	A
C	A	C	A	C	A	C	A	C	A
A	C	A	C	A	C	A	C	A	C
\$	G	\$	G	\$	G	\$	G	\$	G
T	T	T	T	T	T	T	T	T	T
A	T	A	T	A	T	A	T	A	T
A	\$	A	\$	A	\$	A	\$	A	\$

Goal: Find the read “GAT” in the genome “GATTACA”.

Step 3: The two letters that come before “T” in the text are “T” and “A”. Since we’re looking for “GAT”, we’re interested in the “A”.

Finding a word with BWT

B-W	sorted								
T	A	T	A	T	A	T	A	T	A
G	A	G	A	G	A	G	A	G	A
C	A	C	A	C	A	C	A	C	A
A	C	A	C	A	C	A	C	A	C
\$	G								
T	T	T	T	T	T	T	T	T	T
A	T	A	T	A	T	A	T	A	T
A	\$								

Goal: Find the read “GAT” in the genome “GATTACA”.

Step 4: We can repeat the process to find what letter(s) come before “AT” in the text. But first, we need to find which “A” in the sorted BWT corresponds to the “A” in the BWT we found earlier.

Finding a word with BWT

B-W	sorted								
T	A	T	A	T	1 A	T	A	T	A
G	A	G	A	G	2 A	G	A	G	A
C	A	C	A	C	3 A	C	A	C	A
A	C	A	C	1 A	C	A	C	A	C
\$	G								
T	T	T	T	T	T	T	T	T	T
A	T	A	T	2 A	T	A	T	A	T
A	\$	A	\$	3 A	\$	A	\$	A	\$

Goal: Find the read “GAT” in the genome “GATTACA”.

Step 4: The “A” that we found in the BWT is the second in that column, counting downwards. So it should match the second “A” in the sorted BWT.

Finding a word with BWT

B-W	sorted								
T	A	T	A	T	A	T	A	T	A
G	A	G	A	G	A	G	A	G	A
C	A	C	A	C	A	C	A	C	A
A	C	A	C	A	C	A	C	A	C
\$	G								
T	T	T	T	T	T	T	T	T	T
A	T	A	T	A	T	A	T	A	T
A	\$								

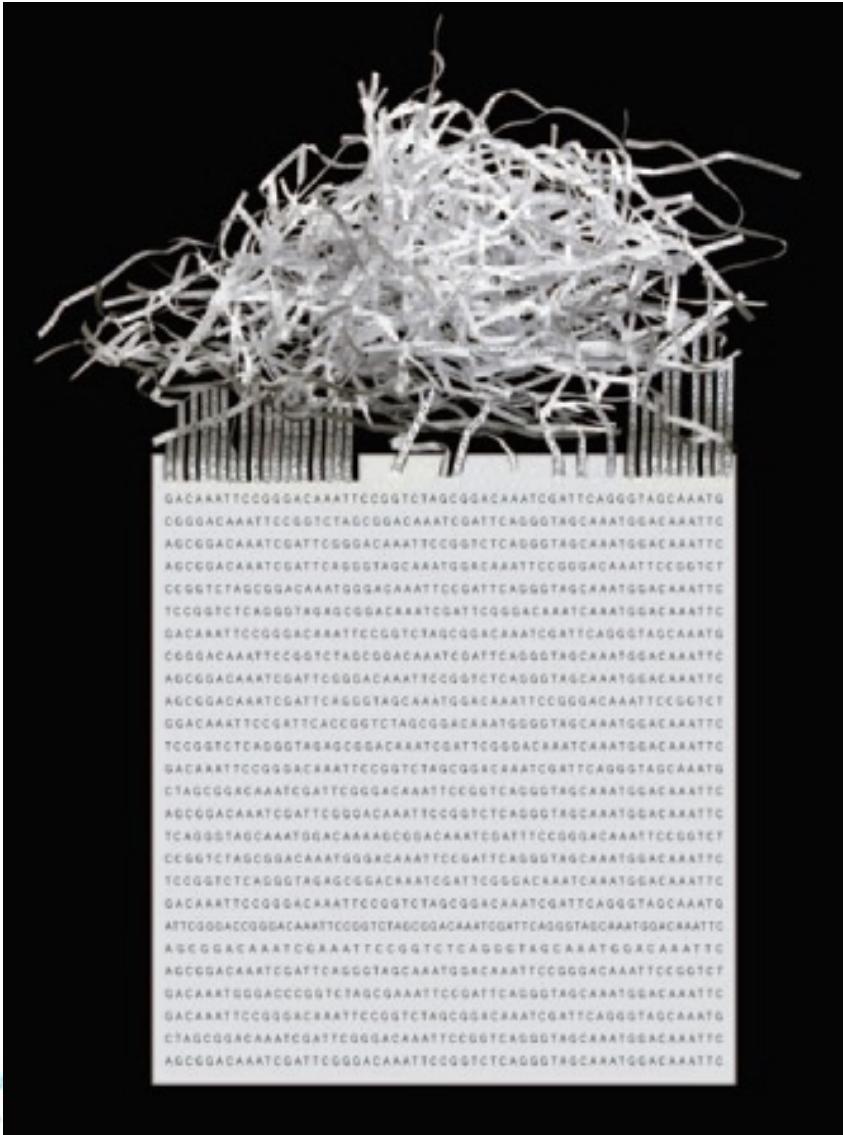
Goal: Find the read “GAT” in the genome “GATTACA”.

Step 5: Repeat Step 3, and find the corresponding letter in the BWT.
Since this is a “G”, we know that “GAT” does exist in our genome.

FM-index

- “Full-text index in Minute space”, invented by Paolo Ferragina and Giovanni Manzini in 2000.
- A BWT-based strategy for quickly searching for patterns within a compressed text.
- Includes additional tricks for finding locations, accounting for mismatches, increasing efficiency.
- FM-indexing lies at the core of popular read mapping software packages such as Bowtie and BWA.

	Bowtie 2	BWA
Index reference genome	<code>bowtie2-build</code>	<code>bwa index</code>
Map reads	<code>bowtie2</code>	<code>bwa mem</code>



Questions?

Source: Kelly Howe, Lawrence Berkeley Laboratory

Baker, Monya. "De novo genome assembly: what every biologist should know." Nature methods 9.4 (2012): 333-337.