





# The biology behind genome assembly

The entire genome cannot be sequenced in one go. Genome assembly (also known as sequence assembly or fragment assembly) refers to the process of aligning and merging fragments from a longer DNA sequence in order to reconstruct the original sequence.

- Whole-genome shotgun sequencing starts by copying and fragmenting (in random order) the DNA followed by sequencing of those fragments, and reassembling them to form a genomic sequence.
- Next-generation sequencing starts by fragmenting DNA/RNA into multiple pieces (non-overlapping), followed by adding adapters, sequencing the libraries in parallel, and reassembling them to form a genomic sequence.

# The biology behind genome assembly

## Whole-genome shotgun sequencing:

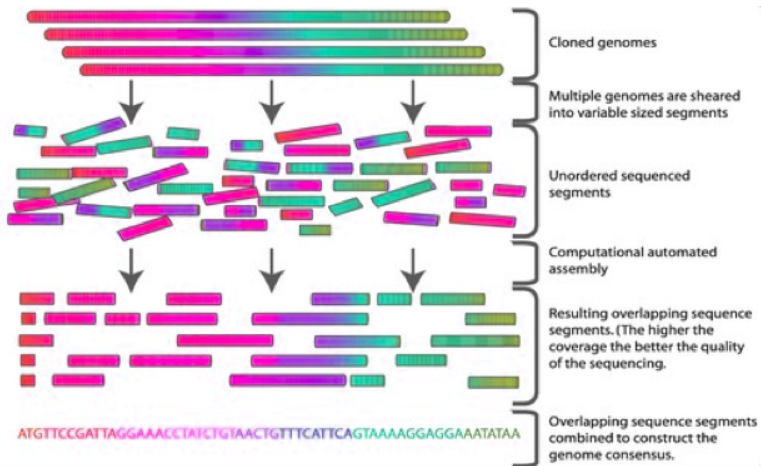
<b>Input:</b>	AAGTCCGTAACCGTCAATTTTCGAGGGCA
<b>Copies:</b>	AAGTCCGTAACCGTCAATTTTCGAGGGCA AAGTCCGTAACCGTCAATTTTCGAGGGCA AAGTCCGTAACCGTCAATTTTCGAGGGCA
<b>Fragments:</b>	AAGT CCGTAACCGT CAATTTTCGA GGGCA AAGTCCG TAACCGTCAA TTTCGAGGGCA AAGTCCG TAACCGTCAATTTTCG AGGGCA

## Next-generation sequencing:

<b>Input:</b>	AAGTCCGTAACCGTCAATTTTCGAGGGCA
<b>Fragments:</b>	AAGT CCGT AACC GTCA ATTT CGAG GGCA



# The mechanism



The entire genome is randomly fragmented and then reassembled

# The coverage

Coverage (also known as read depth or depth) is defined by the average number of reads representing a given character in the consensus sequence. Given the length of the original genome ( $G$ ), number of reads ( $N$ ), and average read length ( $L$ ), one can calculate the coverage ( $\lambda$ ) as follows.

$$\lambda = \frac{N \times L}{G}$$

For example, a hypothetical genome with 10,000 characters (bases) reconstructed from 25 reads with an average length of 800 characters will have 2x redundancy.

# Estimating uncovered bases (Lander-Waterman theory)

Let us assume that the probability that a base is not covered is represented with the following Poisson distribution.

$$\text{Poisson}(0, \lambda) = e^{-\lambda}$$

Then the number of uncovered characters (bases) can be estimated as

$$Ge^{-\lambda}.$$

Similarly, the number of gaps can be estimated as

$$Ne^{-\lambda}$$















# The approach of assembly

The fragment assembly methods for short reads are of two types:

- 1 Overlap Layout Consensus (OLC) assembly
  - Employs string graph-based assemblers
  - Constructs overlap graph from short reads
  - Eliminates redundant reads
  - Traces paths in the graph for assembly
  - Examples include SGA, Fermi, etc.
- 2 de Bruijn Graph-based (DBG) assembly
  - Employs de Bruijn graph-based assemblers
  - Constructs k-mer graph from short reads
  - Discards original reads
  - Traces paths in the graph for assembly

# The approach of assembly

## OLC assembly



Overlap



Layout



Consensus



## DBG assembly



Error correction



de Bruijn Graph



Refine



Scaffolding

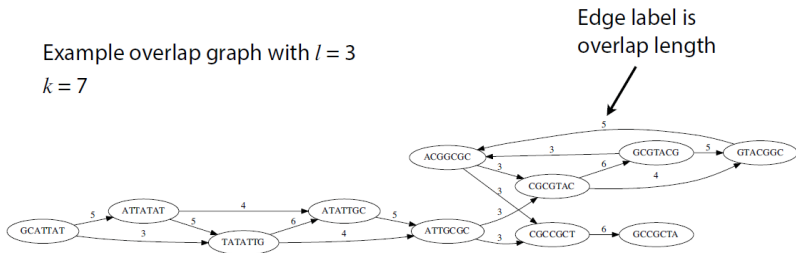


# The OLC assembly

## Step 1 – Overlap: Building an overlap graph

Example overlap graph with  $l = 3$

$k = 7$



**Note:** A merged FM index of all reads enables to match prefixes and suffixes of reads to explore the overlaps.

# The OLC assembly

**Step 2 – Layout:** Putting together stretches of the overlap graph into contigs

We solve the shortest common superstring problem on the overlap graph by revising the graph with edges having a cost of  $-(\text{length of overlap})$ . The shortest common superstring corresponds to a path that visits every node once (Hamiltonian path), minimizing the cost of traversal.

**Note:** Hamiltonian cycle problem is NP-complete.

# The OLC assembly

**Step 3 – Consensus:** Choose the most probable nucleotide sequence for each contig

We remove the transitive edges from the graph.

# The FM-index

An FM-index is a compressed full-text substring index based on the Burrows–Wheeler transform (BWT). The BWT algorithm and its inverse version work as follows.

```
function BWT(string s){
  R <- All possible rotations of s as elements // Right circular shift
  SortedR <- The elements of R in lexicographic order
  return>Last character of each element in SortedR in the order they appear)
}

function InverseBWT(string s){
  SortedT <- {NULL, NULL, ... length(s) times}
  for(i = 1 to length(s)){
    T <- Insert characters of s as prefixes of elements in SortedT
    SortedT <- The elements in T in lexicographic order
  }
  return(The element in SortedT that ends with \$)
}
```

# Burrows–Wheeler Transform (BWT) – An example

^CAGAGAS\$	^CAGAGAS\$	AGAGAS^C	CGG^AA\$A
	\$^CAGAGA	AGAS^CAG	
	A\$^CAGAG	A\$^CAGAG	
	GA\$^CAGA	CAGAGAS^	
	AGAS^CAG	GAGAS^CA	
	GAGAS^CA	GA\$^CAGA	
	AGAGAS^C	^CAGAGAS\$	
	CAGAGAS^	\$^CAGAGA	
<b>Input</b>	<b>Rotate</b>	<b>Sort</b>	<b>Output</b>

**Note:** Two-dimensional array can be used as a data structure.

# Inverse BWT – An example

CGG^AA\$A	C	A	CA	AG	...	AGAGAS^C
	G	A	GA	AG		AGAS^CAG
	G	A	GA	A\$		A\$^CAGAG
	^	C	^C	CA		CAGAGAS^
	A	G	AG	GA		GAGAS^CA
	A	G	AG	GA		GAS^CAGA
	\$	^	\$^	^C		<b>^CAGAGAS</b>
	A	\$	A\$	\$^		\$^CAGAGA
<b>Input</b>	<b>Insert 1</b>	<b>Sort 1</b>	<b>Insert 2</b>	<b>Sort 2</b>	<b>...</b>	<b>Sort 8</b>

# The DBG assembly

**Step 1 – Error Correction:** Correcting errors in the reads

Similar to assembly for whole-genome shotgun sequencing

# The DBG assembly

## Step 2 – de Bruijn Graph: Building a de Bruijn graph

Given a sequence  $s$  and the  $k$ -mer length  $k$ , we can construct a de Bruijn graph with the following characteristics:

- the nodes are  $(k-1)$ -mers, and
- an edge is present between a pair of nodes if they have  $(k-2)$ -long overlap.

Consider the sequence AAATTTA and the value of  $k = 3$ . Let us construct the corresponding de Bruijn graph.

**Note:** The value of  $k$  is taken as odd.

# The DBG assembly

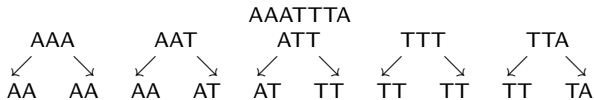
## Step 2 – de Bruijn Graph: Building a de Bruijn graph (continued)

We can obtain the left and right  $(k-1)$ -mers of each  $k$ -mer present in the sequence as follows.

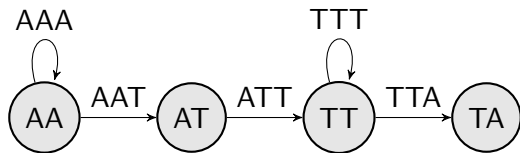
Sequence:

3-mers:

Left/Right 2-mers:



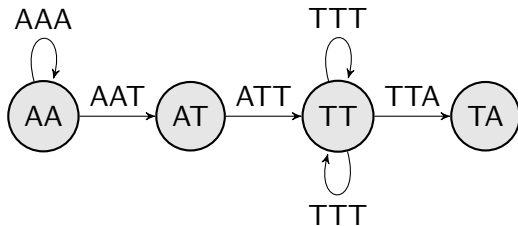
We then construct the corresponding de Bruijn graph as follows.



# The DBG assembly

## Step 2 – de Bruijn Graph: Building a de Bruijn graph (continued)

Note that, if we add one more T to our input sequence such that  $s$  becomes AAATTTTA, and reconstruct the corresponding de Bruijn graph, we get a *multiedge* in the graph as follows.



# The DBG assembly

## Step 2 – de Bruijn Graph: Building a de Bruijn graph (continued)

A graph is *connected* if there is a path between every pair of vertex. An *Eulerian walk* visits each edge exactly once. A directed, connected graph is *Eulerian* if it contains an Eulerian walk.

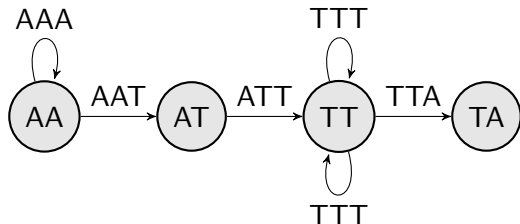
A node is *balanced* if its in-degree = out-degree. A node is *semi-balanced* if its  $|\text{in-degree} - \text{out-degree}| = 1$ . A directed, connected graph is Eulerian if and only if it has at most 2 semi-balanced nodes and all other nodes are balanced.

**Note:** For simplicity, we do not distinguish Eulerian from semi-Eulerian.

# The DBG assembly

## Step 2 – de Bruijn Graph: Building a de Bruijn graph (continued)

Note that, the following graph is Eulerian.



Arguments:

- ①  $AA \rightarrow AA \rightarrow AT \rightarrow TT \rightarrow TT \rightarrow TA$
- ② AA and TA are semi-balanced, AT and TT are balanced

# The DBG assembly

## Step 2 – de Bruijn Graph: Building a de Bruijn graph (continued)

For genome assembly, each  $k$ -mer is recorded in *twin* nodes – one node in forward direction and one node in reverse complement.

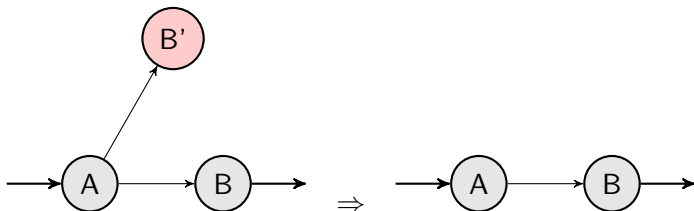
With perfect sequencing, this procedure always yields an Eulerian graph unless the genome is circular. Hence, we just need to find out the Eulerian walk in the graph.

**Note:** No node can be its own reverse complement because  $k$  is odd.



# The DBG assembly

## Step 3 – Refine: Refining the de Bruijn graph

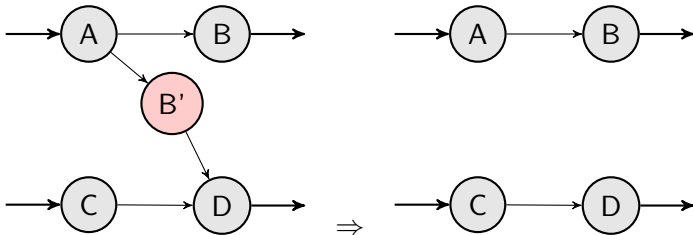


If errors are present at the end of read, remove 'dead-end' tips



# The DBG assembly

## Step 3 – Refine: Refining the de Bruijn graph (continued)



If chimeric edges are present, remove short and low coverage nodes



