

The biology behind genome assembly

Whole-genome shotgun sequencing:

Input:	AAGTCCGTAACCGTCAATTTTCGAGGGCA
Copies:	AAGTCCGTAACCGTCAATTTTCGAGGGCA AAGTCCGTAACCGTCAATTTTCGAGGGCA AAGTCCGTAACCGTCAATTTTCGAGGGCA
Fragments:	AAGT CCGTAACCGT CAATTTTCGA GGGCA AAGTCCG TAACCGTCAA TTTCGAGGGCA AAGTCCG TAACCGTCAATTTTCG AGGGCA

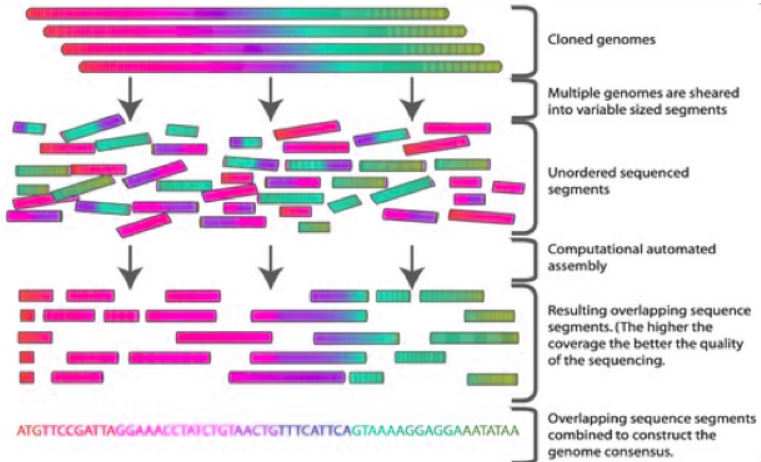
Next-generation sequencing:

Input:	AAGTCCGTAACCGTCAATTTTCGAGGGCA
Fragments:	AAGT CCGT AACC GTCA ATTT CGAG GGCA

Basic terminologies

- **Read:** A 200 to 800 character long word that comes out of a sequencer.
- **Mate pair:** A pair of reads from two ends of the same insert fragment.
- **Contig:** A contiguous sequence formed by several overlapping reads with no gaps.
- **Supercontig (scaffold):** An ordered and oriented set of contigs, usually by mate pairs.
- **Consensus sequence:** A sequence derived from the multiple alignment of reads in a contig.

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 (λ) 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.

The approach of assembly

Even if a pair of reads actually originate from overlapping stretches of the genome, there might be differences.

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

↓

As a result, we have the following issues.

- ① Sequencing error
- ② Difference between inserted copies of chromosome

Note: Humans are diploid (have two copies of each chromosome – one from mother, one from father). The copies can differ.

The approach of assembly

Step 2: Correcting errors (continued)

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

Correlated errors – probably due to repeats



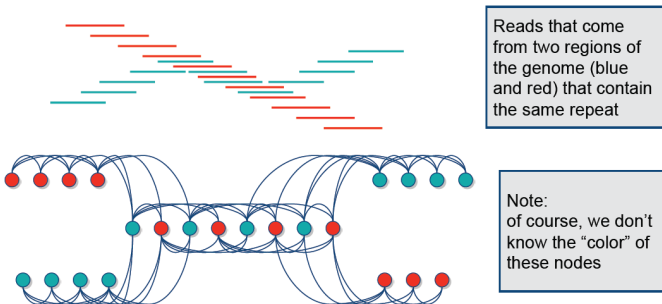
Disentangle overlaps

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

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

The approach of assembly

Step 3: Merging reads



A note

Next-generation sequencing can be of two types as listed below.

- Short-read sequencing (e.g., Illumina sequencing, Roche 454 sequencing)
- Long-read sequencing (e.g., Nanopore sequencing)

The methods of genome assembly are same in general for both the aforementioned types. However, for obvious reasons, assembly of shorter reads are more challenging.

Basic terminologies

- **Read:** A 50 to 300 character long (for short-read sequencing) or 500 to 2.3 million character long (for long-read sequencing) word that comes out of a sequencer.
- **Contig:** A contiguous sequence formed by several overlapping reads with no gaps.
- **Scaffold:** An ordered and oriented set of contigs.
- **Consensus sequence:** A sequence derived from the multiple alignment of reads.

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

An overlap graph is a directed graph which is constructed with the following considerations.

- A node \equiv A read
- A edge \equiv An overlap of reads

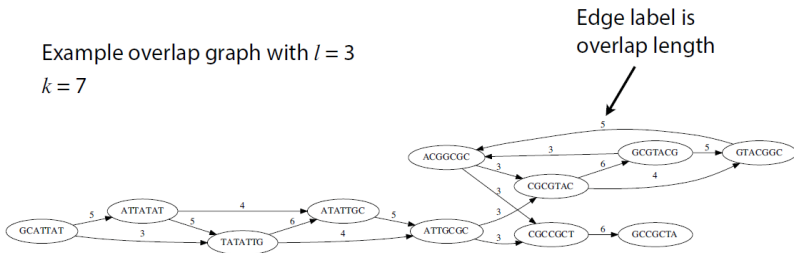
An edge (r_1, r_2) signifies that a suffix of the read r_1 (source) is similar to a prefix of the read r_2 (sink).

Note: An overlap graph may contain a cycle.

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 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

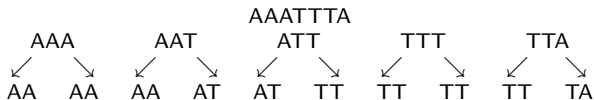
Step 2 – de Bruijn Graph: Building a de Bruijn graph (continued)

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

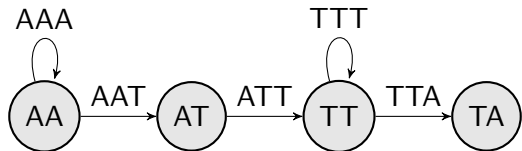
Sequence:

3-mers:

Left/Right 2-mers:



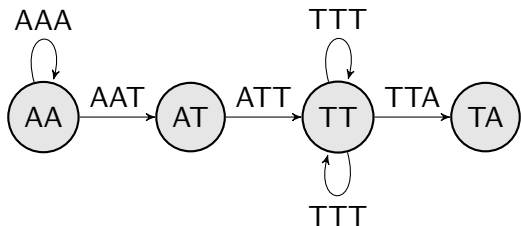
We then construct the corresponding de Bruijn graph as follows.



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.

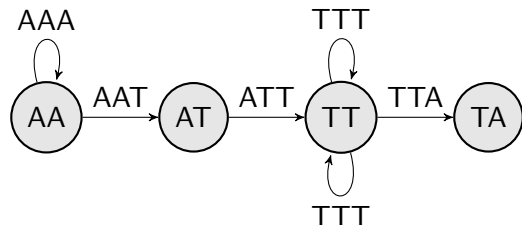
A node is *balanced* if its in-degree = out-degree. A node is *semi-balanced* if its $|\text{in-degree} - \text{out-degree}| = 1$. A directed, connected graph is Eulerian if and only if it has at most 2 semi-balanced nodes and all other nodes are balanced.

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

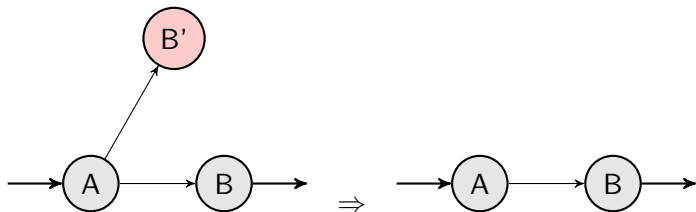
Step 2 – de Bruijn Graph: Building a de Bruijn graph (continued)

Note that, there are further challenges to deal with as listed below.

- 1 A de Bruijn graph can have multiple Eulerian walks, *only one* of which corresponds to original superstring.
- 2 Having gaps in the coverage can lead to a disconnected de Bruijn graph. The connected components are individually Eulerian, but not the overall graph.
- 3 Differences in coverage might lead to non-Eulerian de Bruijn graphs.
- 4 Errors and differences between chromosomes might lead to non-Eulerian de Bruijn graphs.

The DBG assembly

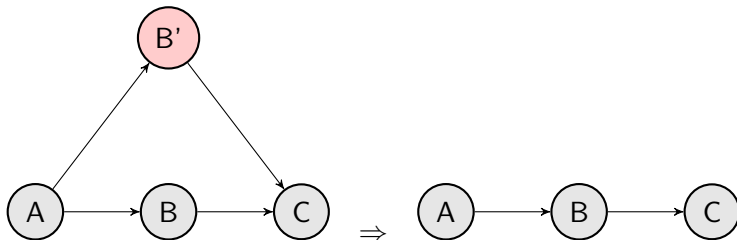
Step 3 – Refine: Refining the de Bruijn graph



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

The DBG assembly

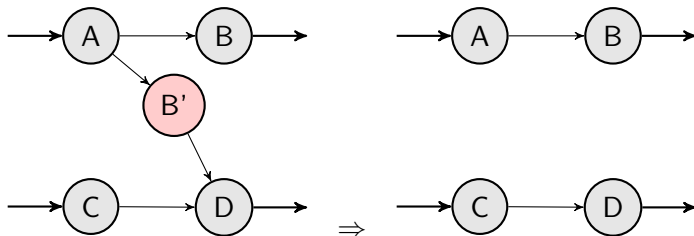
Step 3 – Refine: Refining the de Bruijn graph (continued)



If errors are present in the middle of read, remove 'bubbles'

The DBG assembly

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



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

