# Introduction to genome assembly

#### CMSC423

Many slides courtesy of Ben Langmead and Carl Kingsford

### SEC-GEN SEQUENCING





### SEC-GEN SEQUENCING



Fragmentation is random, i.e., not equal-sized (but hard to draw)



#### SEC-GEN SEQUENCING





# SEQUENCING

- "Ultra high throughput" DNA sequencing
  - 6 gigabases / day vs.
  - 3 gigabases / 13 years (human genome project, more or less)
  - 200 bp long reads



## From reads to evidence

#### 0 🔴 💮

@HWI-EAS146:5:1:1:961#0/1 @HWI-EAS146:5:1:1:1595#0/1 @HWI-EAS146:5:1:1:1048#0/1 CTGGACTGCATCCTACCACCAACTCGTCCAANNNNCNNNNNNNNN @HWI-EAS146:5:1:1:1607#0/1 CTCCTCTCAAGGTCCCCAGAAGCACAGCCAANNNNANTNNCTNNNN @HWI-EAS146:5:1:1:1719#0/1 CACGATCTGGGTTTATTGTAACCTCCGCCTCNNNNGNTNAAGNNNN @HWI-EAS146:5:1:2:947#0/1 CCCAGGAGAAAGCCATGTTCAGTTCGAGCGCNNANANCGTGANNNN @HWI-EAS146:5:1:2:563#0/1 CCAGCCCCTCCCCCATCTCCCACCCTGTACCTNANCCCCTGANNNN @HWI-EAS146:5:1:2:1631#0/1 TGGGAACGCAGCCTACACTCTTCCCAGGCCTCCTNCCTCCGTNNNN BBB@6@BBBBBBBB@BBBABAABBB?:9BB@BA5&<B:%%%%%% @HWI-EAS146:5:1:2:1420#0/1 CTCAAACTCCTGACCTTTGGTGATCCACCCGCCTNGGCCTTCNNNN @HWI-EAS146:5:1:1:961#0/1 @HWI-EAS146:5:1:1:1595#0/1 @HWI-EAS146:5:1:1:1048#0/1 @HWI-EAS146:5:1:1:1607#0/1 CTCCTCTCAAGGTCCCCAGAAGCACAGCCAANNNNANTNNCTNNNN

+ BBCCCCCCBBCB7CBC=7>+<>=BCBCB%%%%%%%%%%%%%%%%

## From reads to evidence

#### 0 0

+

+ BCC?+<B=?BB5=ABA?B6BBBB4BB?B%%%%%%%%%%%%%% @HWI-EAS146:5:1:2:947#0/1

CCCAGGAGAAAGCCATGTTCAGTTCGAGCGCNNANANCGTGANNN +

CCAGCCCCTCCCCCATCTCCCACCCTGTACCTNANCCCCTGANNNN +

BBB@6@BBBBBBBBB@BBBABAABBB?;9BB@BA5&<B:%%%%%% @HWI-EAS146:5:1:2:1420#0/1 CTCAAACTCCTGACCTTTGGTGATCCACCCGCCTNGGCCTTCNNNN

+ BBBB:BBBBBABAAA?:(=A8@>AAA?AB?=A%%%%%%%%%%%%%%% @HWI-EAS146:5:1:1:961#0/1

CTCCTCTCAAGGTCCCCAGAAGCACAGCCAANNNNANTNNCTNNNN + BBCCCCCCBBCB7CBC=7>+<>=BCBCB%%%%%%%%%%%%%%%%%

@HWI-EAS146:5:1:1:1719#0/1 CACGATCTGGGTTTATTGTAACCTCCGCCTCNNNNGNTNAAGNNNN

#### I. de novo

Assume nothing! - let reads tell us everything



Source: De Novo Assembly Using Illumina Reads. Illumina. 2010 http://www.illumina.com/Documents/products/technotes/technote\_denovo\_assembly.pdf

# What we'll cover

- Genome assembly as graph problems
  - Two representations:
  - Overlap graph
    - How much sequencing required for assembly
  - DeBruijn graph
- How to get assemblies from solutions to graph problems

#### Overlap Graph

Overlap graph: Nodes = reads Edges = overlaps



Given overlap graph, how can we find a good candidate assembly?

#### Overlap Graph



Given overlap graph, how can we find a good candidate assembly?

#### Overlap Graph



Given overlap graph, how can we find a good candidate assembly?

Hamiltonian Path (aka Traveling Salesman Path): visit every node in the graph exactly once.

### Hamiltonian Path

- Motivation: Every read must be used in exactly one place in the genome.
- Hamiltonian Path is NP-hard.
- Though good solvers exist, they can't operate on the millions of reads from a sequencing project.
- Solution: greedy walk along the graph.



#### Shotgun Sequencing

Many copies of the DNA



Shear it, randomly breaking them into many small pieces, read ends of each:

Assemble into original genome:

#### Lander-Waterman Statistics

How many reads to we need to be sure we cover the whole genome?



An *island* is a contiguous group of reads that are connected by overlaps of length  $\geq \theta L$ . (Various colors above)

Want: Expression for expected # of islands given  $N, g, L, \theta$ .

#### Expected # of Islands

 $\lambda := N/g =$  probability a read starts at a given position (assuming random sampling)

#### Pr(*k* reads start in an interval of length *x*)

*x* trials, want *k* "successes," small probability  $\lambda$  of success Expected # of successes =  $\lambda x$ Poisson approximation to binomial distribution:

$$\Pr(k \text{ reads in length } x) = e^{-\lambda x} \frac{(\lambda x)^k}{k!}$$

Expected # of islands =  $N \times Pr(read is at rightmost end of island)$ 

$$(1-\theta)L \quad \theta L = N \times \Pr(\text{o reads start in } (1-\theta)L)$$
$$= Ne^{-\lambda(1-\theta)L} \frac{(\lambda(1-\theta)L)^{0}}{0!}$$
$$= Ne^{-\lambda(1-\theta)L}$$
$$= Ne^{-(1-\theta)LN/g} \leftarrow LN/g \text{ is called the coverage } c.$$

#### Expected # of Islands, 2

Rewrite to depend more directly on the things we can control: c and  $\theta$ 

Expected # of islands =  $Ne^{-(1-\theta)LN/g}$ 

 $= N e^{-(1-\theta)c}$ 



#### Assembly via <u>Eulerian</u> Path

## de Bruijn graph



### Example bacterial de Bruijn graph



Paths with no branches compressed into a single node

#### Eulerian path =

use every edge exactly once.

With perfect data, the genome can be reconstructed by some Eulerian path through this graph

#### Assembly via Eulerian Path



Let dG(s) be the de Bruijn graph of string s. Then s corresponds to some Eulerian path in dG(s).

A directed graph has an Eulerian path if and only if:

- •One node has one more edge leaving it than entering
- •One node has one more edge entering than leaving
- •All other nodes have the same number of edges entering and leaving

How can we find such a path?



#### A directed graph has an Eulerian cycle if and only if:

•All nodes have the same number of edges entering and leaving

tagacgaacgtacggtagg



## Eulerian Path Algorithm

Connect node with out-degree < in-degree to node with out-degree < in-degree. So that we will have an Eulerian cycle.

Why will you return to *u*?

Walk from some arbitrary node *u* until you return to *u*, creating a doubly liked list of the path you visit.

#### **Repeat** until all edges used:

\*How can find such a node quickly?

•Start from some node *w* on the current tour with unused edges<sup>\*</sup>.

•Walk along unused edges until you return to w, inserting the visited nodes after w into the current tour list.





## Eulerian Path Algorithm

Connect node with out-degree < in-degree to node with out-degree < in-degree. So that we will have an Eulerian cycle.

Why will you return to *u*?

Walk from some arbitrary node *u* until you return to *u*, creating a doubly liked list of the path you visit.

#### **Repeat** until all edges used:

\*How can find such a node quickly?

•Start from some node *w* on the current tour with unused edges<sup>\*</sup>.

•Walk along unused edges until you return to w, inserting the visited nodes after w into the current tour list.





### The Problem with Eulerian Paths

There are typically an astronomical number of possible Eulerian tours with perfect data.

Adding back constraints to limit # of tours leads to a NP-hard problem.

With imperfect data, there are usually NO Eulerian tours.



Aside: counting # of Eulerian tours in a directed graph is easy, but in an undirected graph is #P-complete (hard).

#### Mate Pairs





Gapped Genome Path String Problem

Reconstruct a string from a sequence of (k,d)-mers corresponding to a path in a paired de Bruijn graph.

Given: A sequence of (k, d)-mers  $(a_1 lb_1), ..., (a_n lb_n)$  such that Suffix $(a_i lb_i) = Prefix(a_{i+1} lb_{i+1})$  for all i from 1 to n-1.

Return: A string Text where the i-th k-mer in Text is equal to  $Suffix(a_i|b_i)$  for all i from 1 to n, if such a string exists.

# References

- <u>http://www.cbcb.umd.edu/research/assembly\_primer</u>
- <u>http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2874646/</u>
- <u>http://www.math.ucsd.edu/~gptesler/186/slides/shotgun\_f13-handout.pdf</u>
- http://www.biomedcentral.com/1471-2105/11/21/abstract