

Genome Assembly

Lecture 20
Oct 19, 2016

Announcements

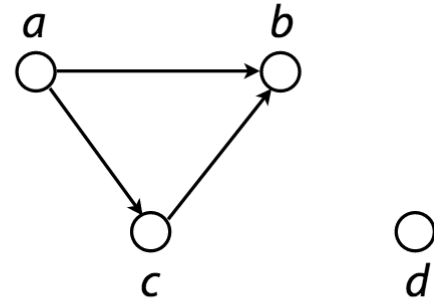
ASSEMBLY

Directed graph $G(V, E)$ consists of set of *vertices*, V and set of *directed edges*, E

Directed edge is an *ordered pair* of vertices.
First is the *source*, second is the *sink*.

Vertex is drawn as a circle

Edge is drawn as a line with an arrow
connecting two circles



Vertex also called *node* or *point*

Edge also called *arc* or *line*

Directed graph also called *digraph*

$$V = \{a, b, c, d\}$$

$$E = \{(a, b), (a, c), (c, b)\}$$

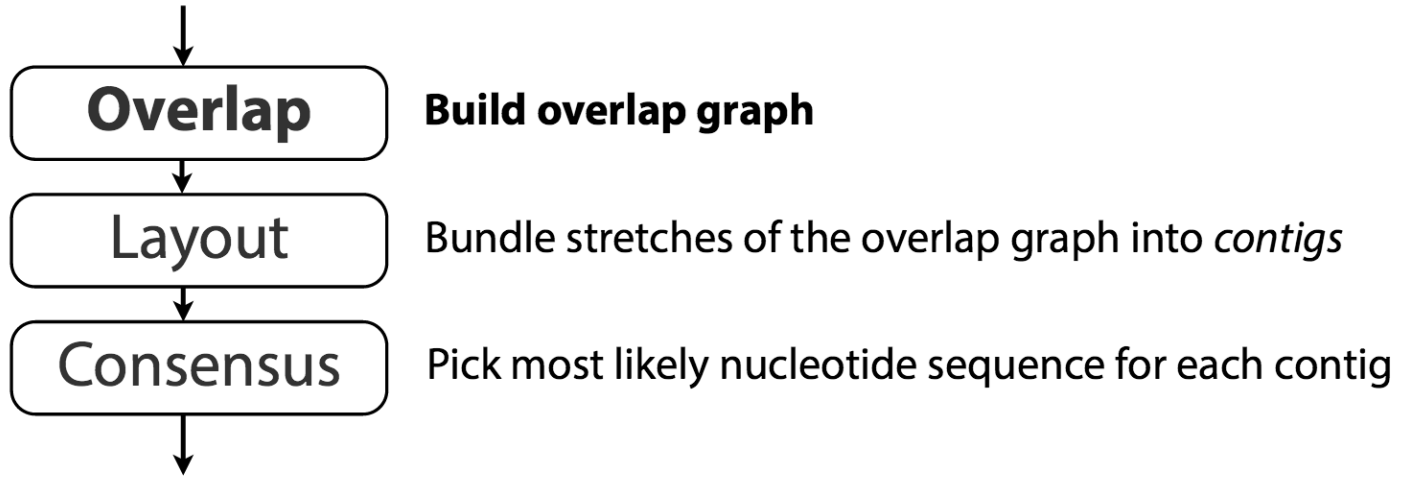
Source Sink

ASSEMBLY

- 2 assembly strategies:

ASSEMBLY

- OLC Assembly



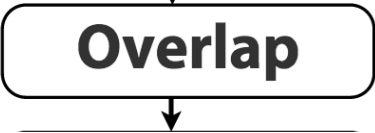
ASSEMBLY

- OLC Assembly: Characteristics

ASSEMBLY

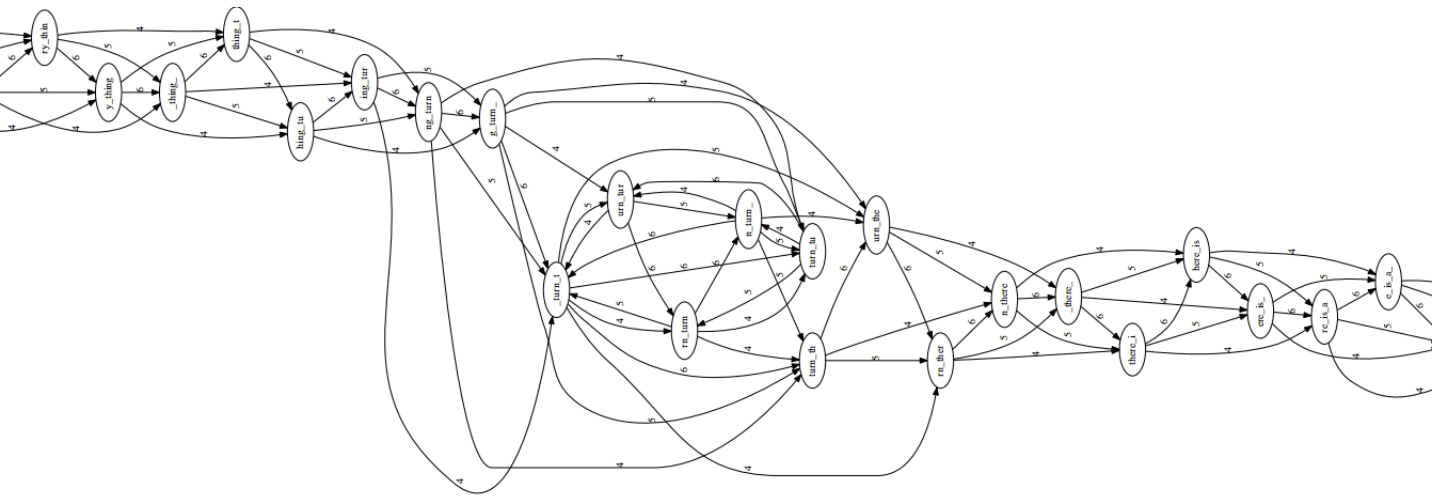
<https://youtu.be/yPJ7yHRk2OI>

ASSEMBLY



Build overlap graph

to_every_thing_turn_turn_turn_there_is_a_season
L=4, k=7



ASSEMBLY

Overlap

Build overlap graph

Vertices (reads): { *a*: CTCTAGGCC, *b*: GCCCTCAAT, *c*: CAATTTTT }

Edges (overlaps): { (*a*, *b*), (*b*, *c*) }

a: CTCTAGGCC

3

b: GCCCTCAAT

4

c: CAATTTTT

CTCTAGGCC

|||

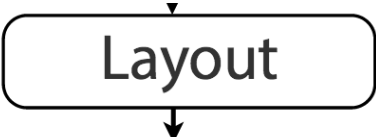
GCCCTCAAT

GCCCTCAAT

||||

CAATTTTT

ASSEMBLY - OLC

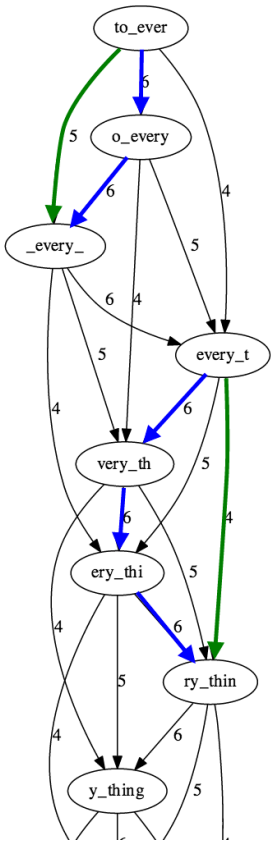


Bundle stretches of the overlap graph into *contigs*

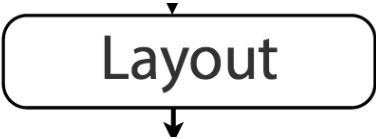
Anything redundant about this part of the overlap graph?

Some edges can be *inferred (transitively)* from other edges

E.g. **green** edge can be inferred from **blue**

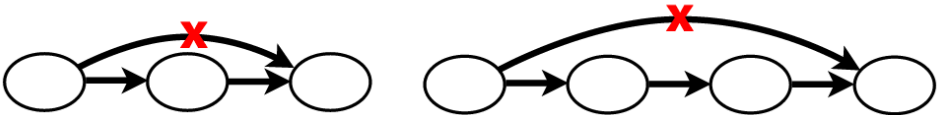


ASSEMBLY - OLC

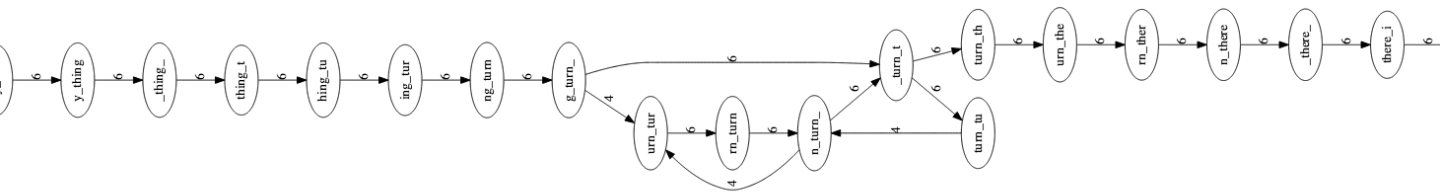


Bundle stretches of the overlap graph into *contigs*

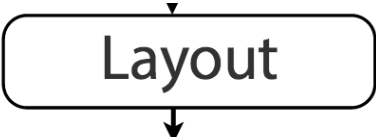
Remove transitively-inferrible edges, starting with edges that skip one or two nodes:



After:

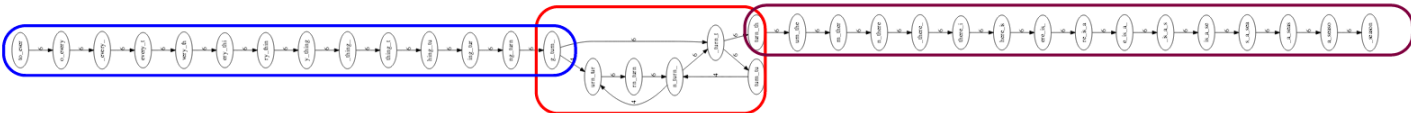


ASSEMBLY - OLC



Bundle stretches of the overlap graph into *contigs*

Emit *contigs* corresponding to the non-branching stretches



Contig 1
to_every_thing_turn_

Contig 2
turn_there_is_a_season

Unresolvable repeat

ASSEMBLY - OLC

Consensus

Pick most likely nucleotide sequence for each contig

TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAACTA
TAG TTACACAGATTATTTGACTTCATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

Take reads that make
up a contig and line
them up

↓ ↓ ↓ ↓ ↓
TAGATTACACAGATTACTGACTTGATGGCGTAA CTA

Take *consensus*, i.e.
majority vote

At each position, ask: what nucleotide (and/or gap) is here?

Complications: (a) sequencing error, (b) ploidy

Say the true genotype is AG, but we have a high sequencing error rate
and only about 6 reads covering the position.

ASSEMBLY – DE BRUIJN

<https://youtu.be/9O3hAXp8gdM?list=PLQ-85IQIPqFNGdaeGpV8dPEeSm3AChb6L>

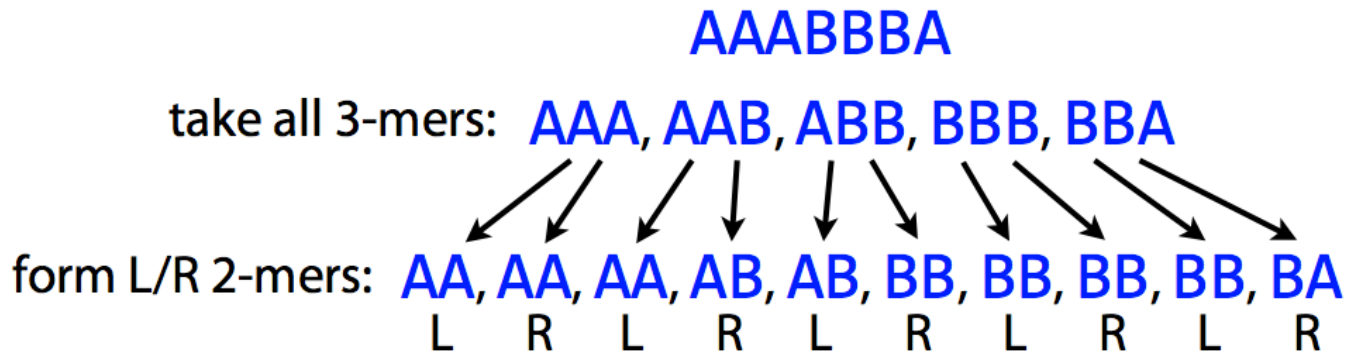
ASSEMBLY – DE BRUIJN

Hamiltonian Path Problem

Eulerian Path Problem

ASSEMBLY – DE BRUIJN

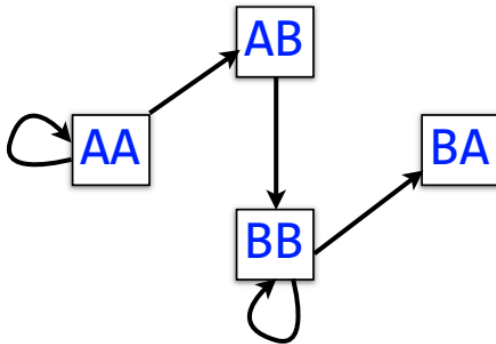
ASSEMBLY – DE BRUIJN



ASSEMBLY – DE BRUIJN

form L/R 2-mers: AA, AA, AA, AB, AB, BB, BB, BB, BB, BA
L R L R L R L R L R

Let 2-mers be nodes in a new graph. Draw a directed edge from each left 2-mer to corresponding right 2-mer:



Each *edge* in this graph corresponds to a length-3 input string