



CS 68: BIOINFORMATICS

Prof. Sara Mathieson
Swarthmore College
Spring 2018



Outline: Jan 26

- Continue: overlap graph assembly
- Begin: de Bruijn graph assembly

Notes:

- Lab 1 due Wednesday
- Assembly reading posted: spend 1.5 hours max

Steps of Overlap Graph Assembly (also called “overlap-layout-consensus”)

- 1) Compute **overlaps between all pairs of reads**. With n = number of reads and L = length of reads, this is naively $O(n^2L^2)$. We will learn better ways of “aligning” sequences next week.
- 2) Construct a **graph with reads as the nodes** and **directed, weighted edges** between reads with $\geq T$ overlap.

Activity example: $L = 10, T = 5$

ATATATACTGGCGTATCGCAGTAAACGCGCCG

R1 : ACTGGCGTAT

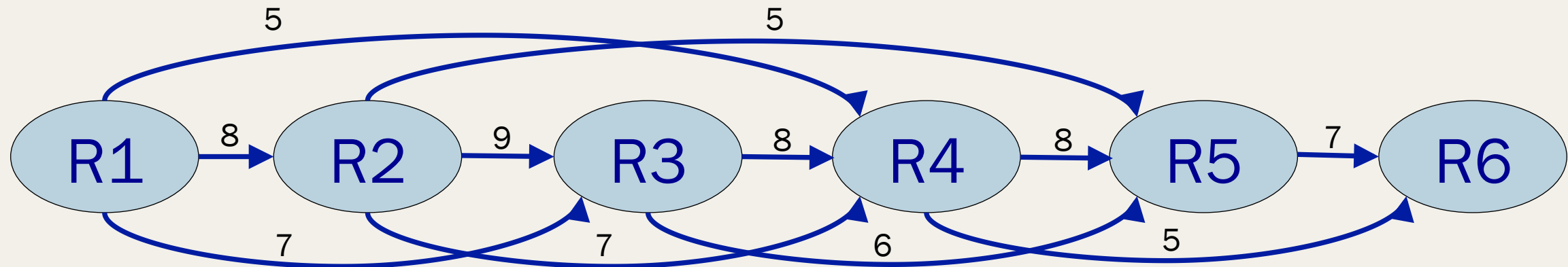
R2 : TGGCGTATCG

R3 : GGCGTATCGC

R4 : CGTATCGCAG

R5 : TATCGCAGTA

R6 : CGCAGTAAAC



Issues with overlap graphs

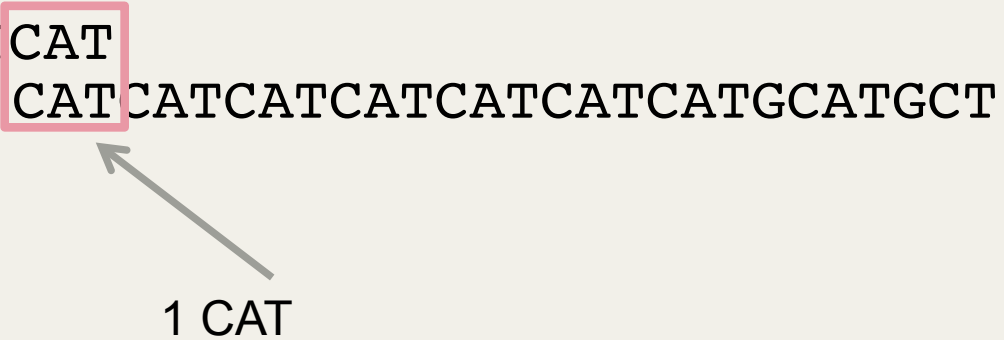
Bubbles



Repeats

read 1: TAACTGTTTCGCATCATCATCAT
read 2: CATCATCATCATCATCATGCATGCT

TAACTGTTTCGCATCATCAT
CATCATCATCATCATCATGCATGCT

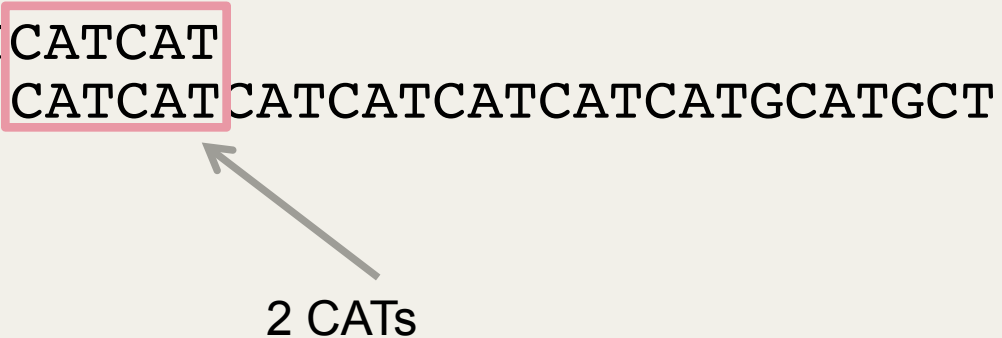


1 CAT

Repeats

read 1: TAACTGTTTCGCATCATCATCAT
read 2: CATCATCATCATCATCATGCATGCT

TAACTGTTTCGCATCATCATCATCATCATCATGCATGCT

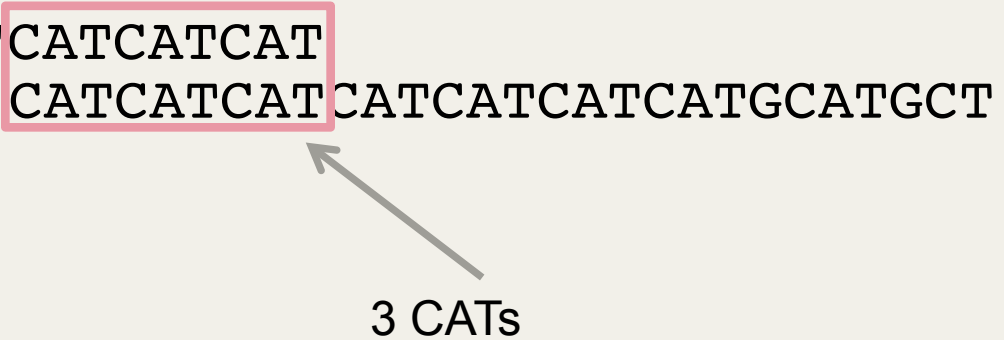


2 CATs

Repeats

read 1: TAACTGTTTCGCATCATCATCAT
read 2: CATCATCATCATCATCATGCATGCT

TAACTGTTTCGCATCATCATCAT
CATCATCATCATCATCATGCATGCT



3 CATs

Repeats

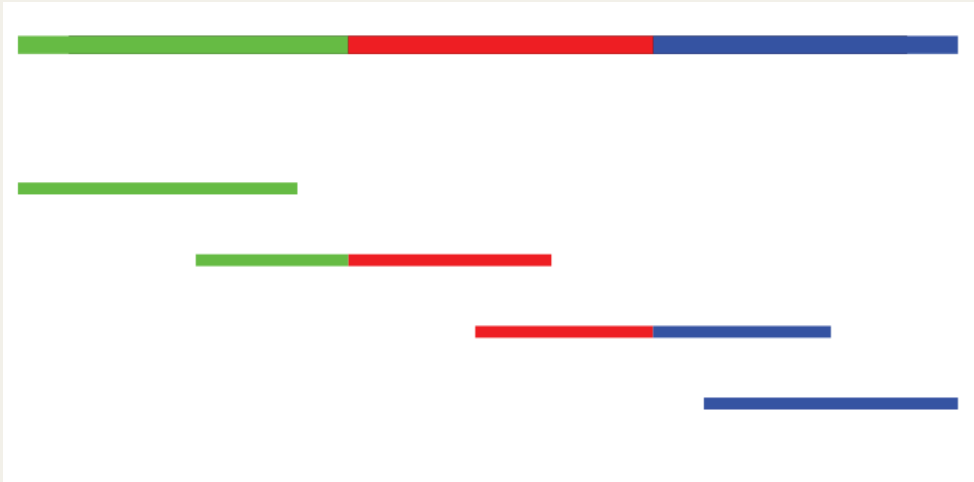
read 1: TAACTGTTTCGCATCATCATCAT
read 2: CATCATCATCATCATCATGCATGCT

TAACTGTTTCGCATCATCATCAT
CATCATCATCATCATCATGCATGCT

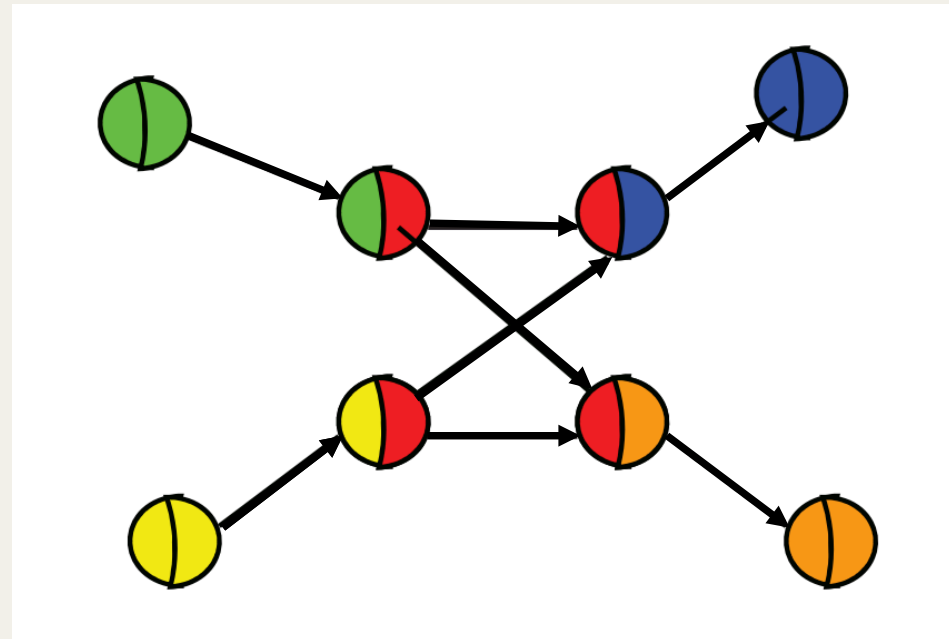
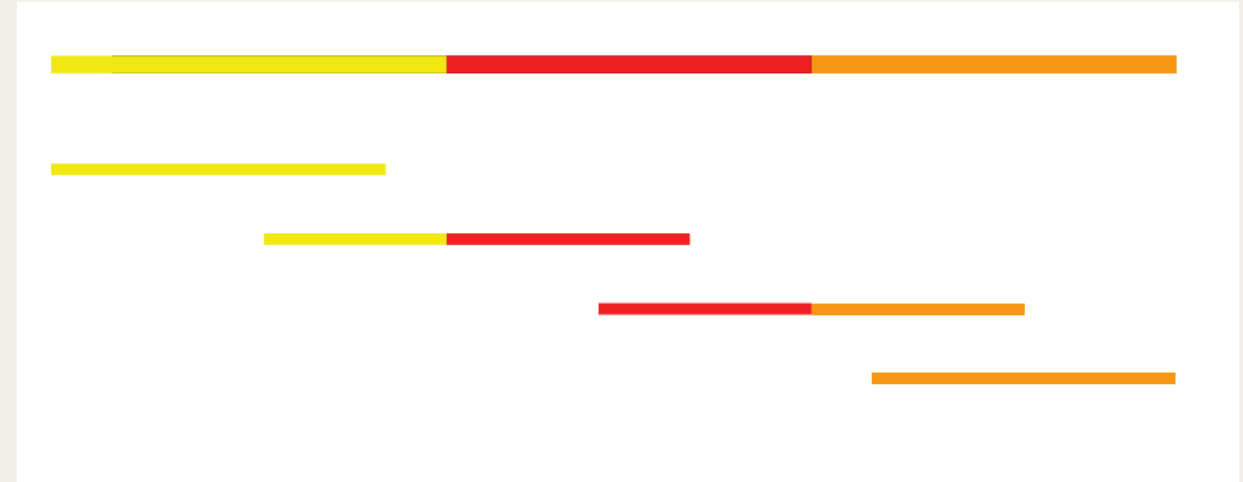
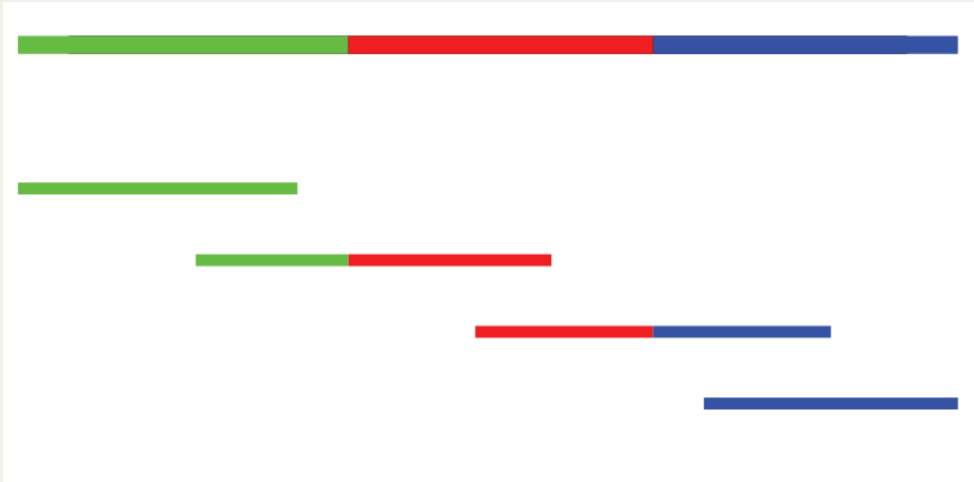


4 CATs

What would the graph look like for these reads?



What would the graph look like for these reads?



Back to overlap graph algorithm

Steps of Overlap Graph Assembly (also called “overlap-layout-consensus”)

- 1) Compute **overlaps between all pairs of reads**. With n = number of reads and L = length of reads, this is naively $O(n^2L^2)$. We will learn better ways of “aligning” sequences next week.
- 2) Construct a **graph with reads as the nodes** and **directed, weighted edges** between reads with $\geq T$ overlap.
- 3) “Layout” the graph and try to “group” stretches of the graph into “**contigs**” (short for contiguous), these are (hopefully) long portions of the original genome
- 4) Find a “consensus” *sequence* for each contig

Activity example: $L = 10, T = 5$

ATATATACTGGCGTATCGCAGTAAACGCGCCG

R1 : ACTGGCGTAT

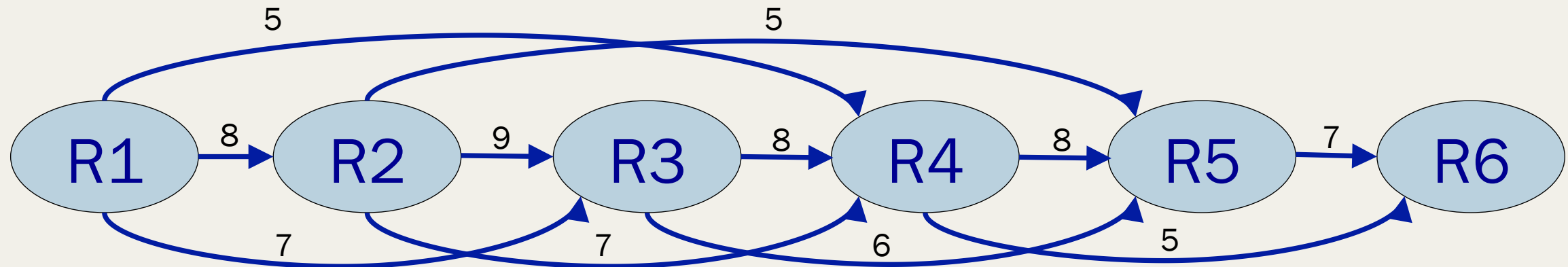
R2 : TGGCGTATCG

R3 : GGCGTATCGC

R4 : CGTATCGCAG

R5 : TATCGCAGTA

R6 : CGCAGTAAAC



Activity example: $L = 10, T = 5$

ATATATACTGGCGTATCGCAGTAAACGCGCCG

R1 : ACTGGCGTAT

R2 : TGGCGTATCG

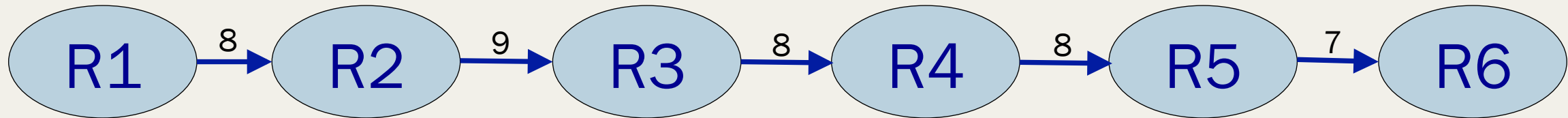
R3 : GGCGTATCGC

R4 : CGTATCGCAG

R5 : TATCGCAGTA

R6 : CGCAGTAAAC

First simplification: remove
edges that can be (transitively)
inferred from other edges



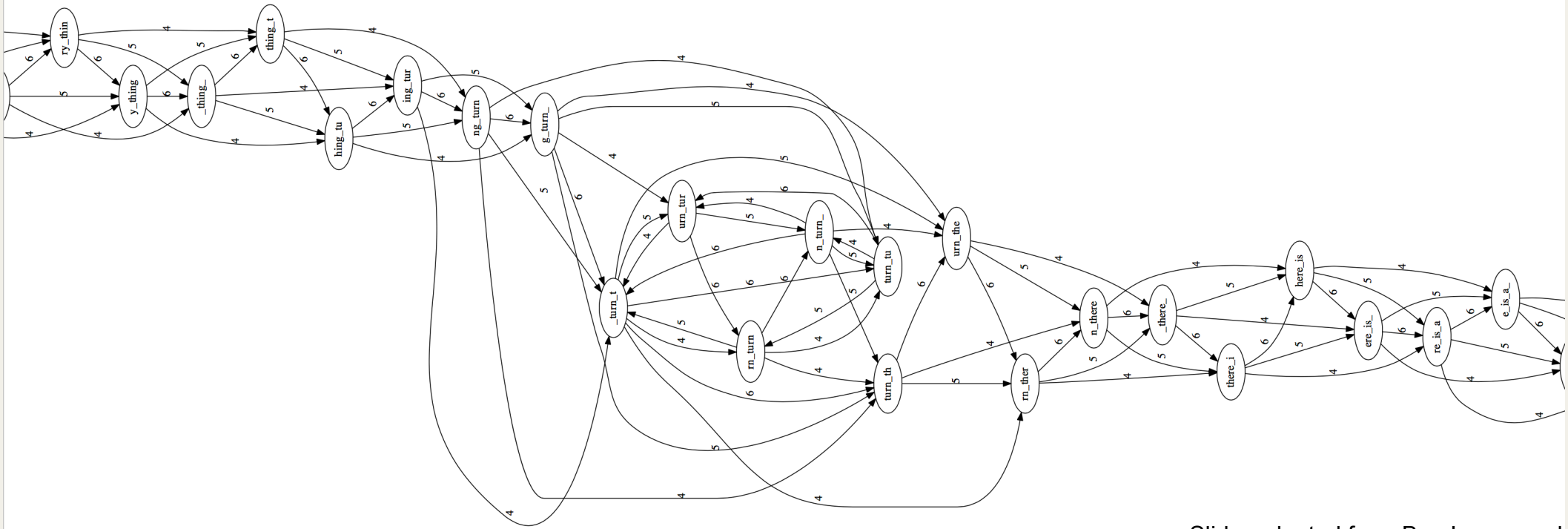
Layout

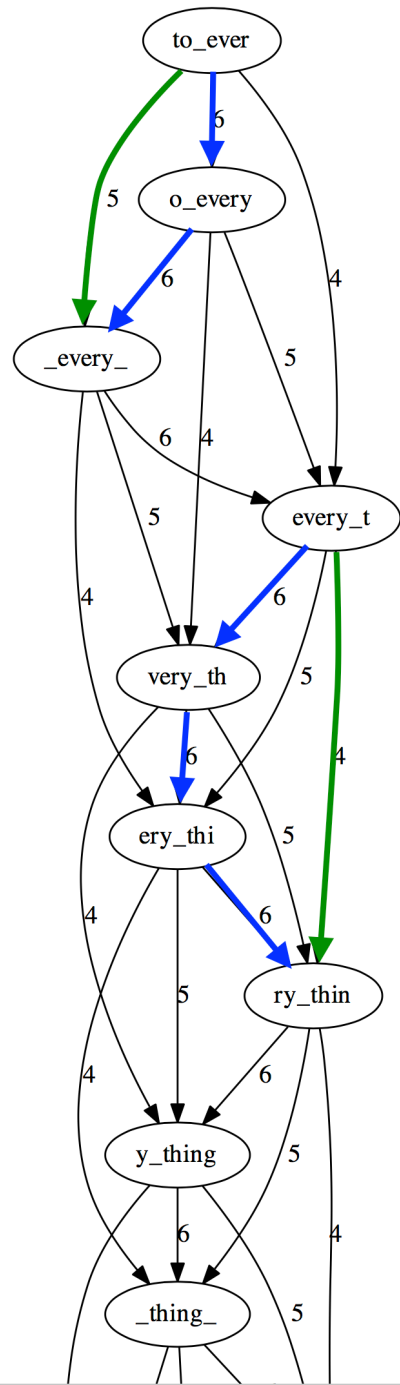
Overlap graph is big and messy. Contigs don't "pop out" at us.

Below: part of the overlap graph for

to_everything_turn_turn_turn_there_is_a_season

$L = 7, T = 4$

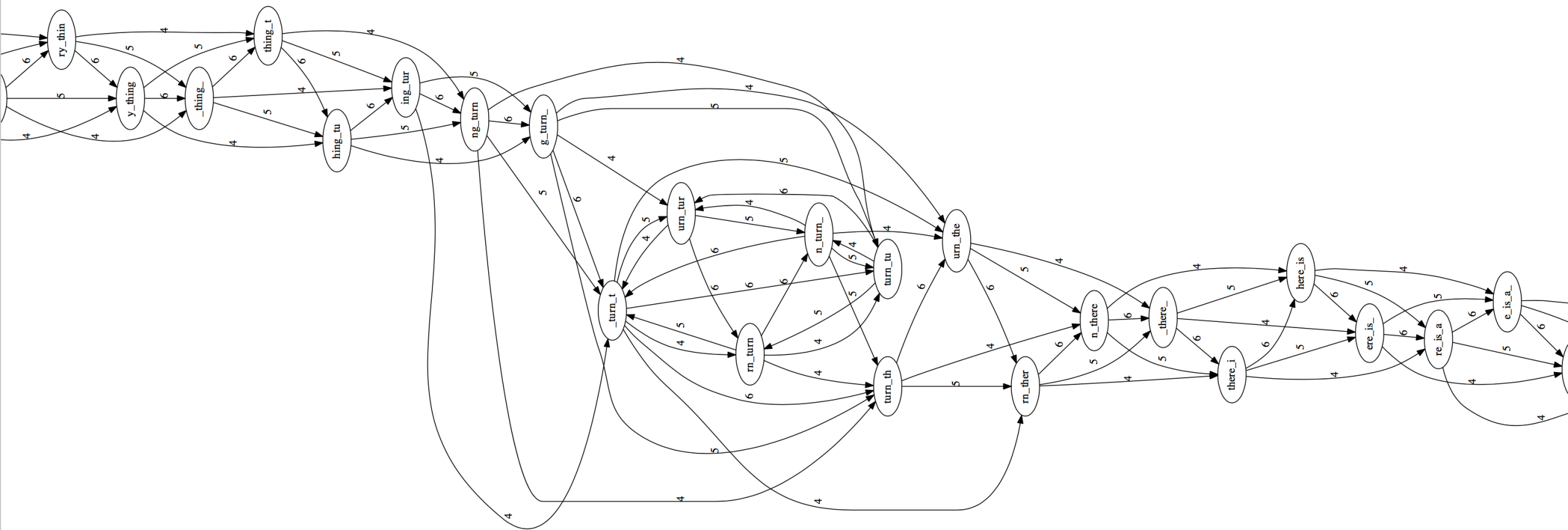




In this example: green edges
can be inferred from blue

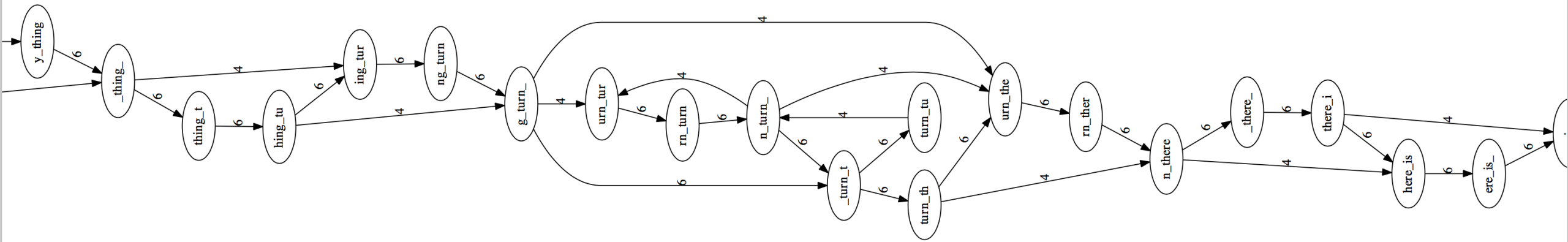
Layout: remove transitively-inferrable edges

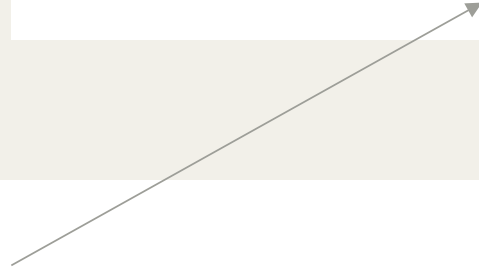
Before:



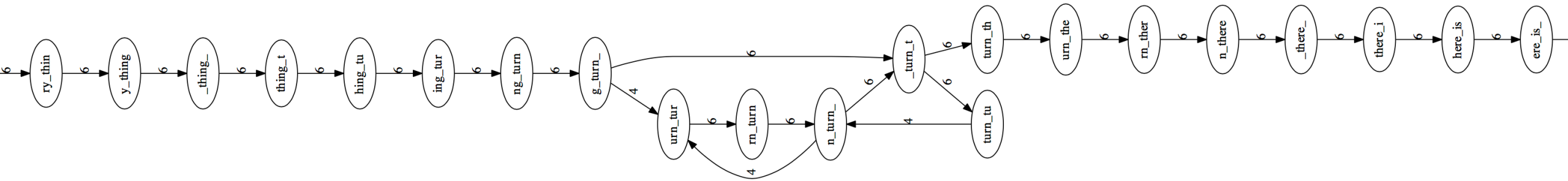
Layout: remove transitively-inferrable edges

After removing edges that skip one node



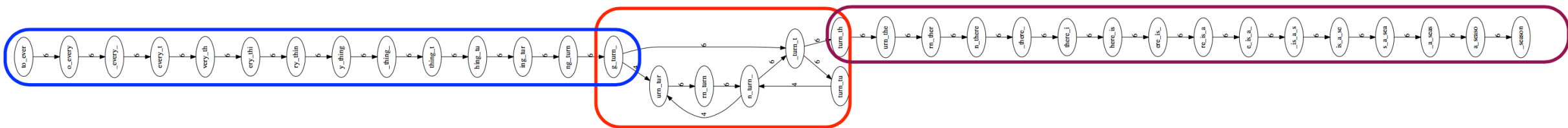


After removing edges that skip one or two nodes



Layout

Emit *contigs* corresponding to the non-branching stretches



Contig 1

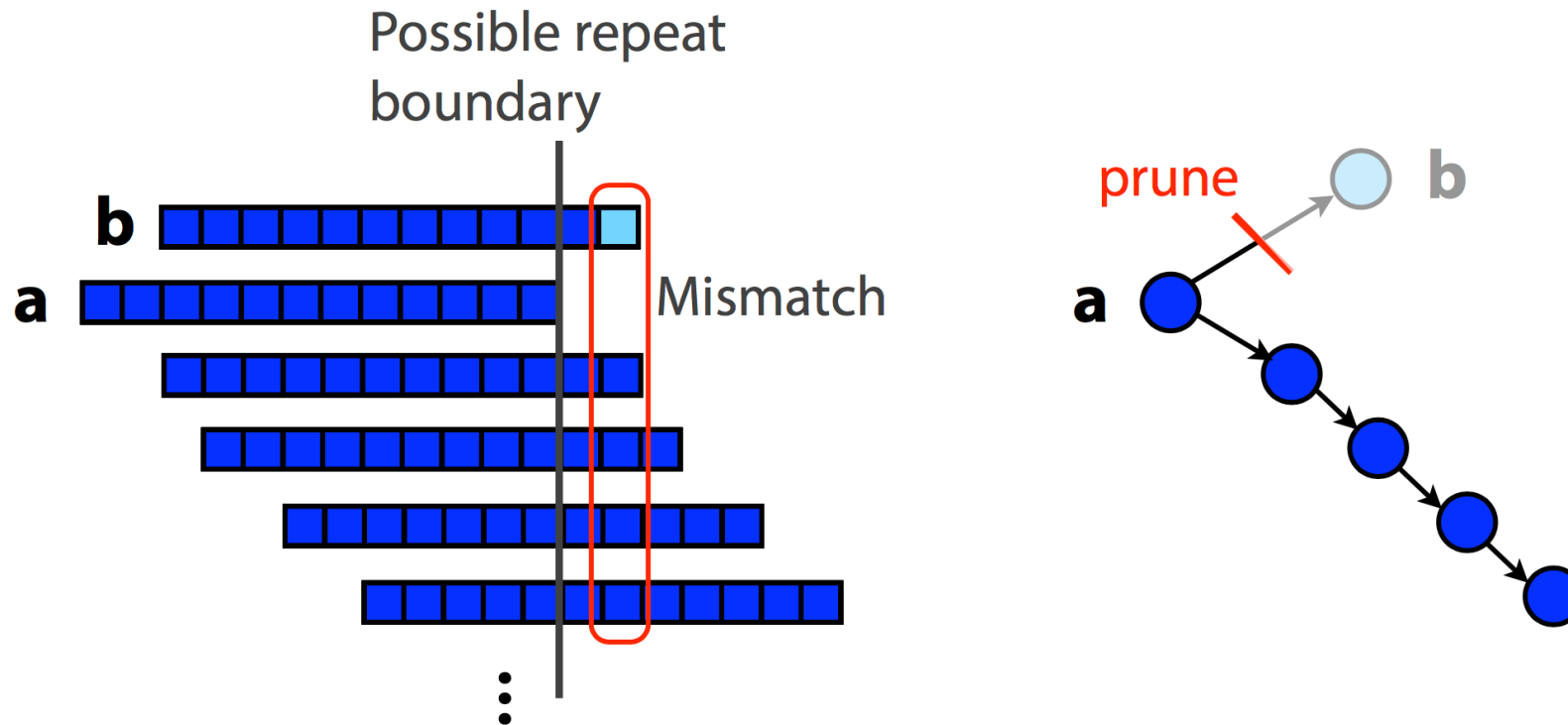
to_every_thing_turn_

Contig 2

turn_there_is_a_season

Unresolvable repeat

In practice, layout step also has to deal with spurious subgraphs, e.g. because of sequencing error



Mismatch could be due to sequencing error or repeat. Since the path through **b** ends abruptly we might conclude it's an error and prune **b**.

Consensus

TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTTGATGGCGTAAACTA
TAG TTACACAGATTATTGACTTTCATGGCGTAA CTA
TAGATTACACAGATTACTGACTTTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTTGATGGCGTAA CTA

↓ ↓ ↓ ↓ ↓
TAGATTACACAGATTACTGACTTTGATGGCGTAA CTA



Take reads that make
up a contig and line
them up

Take *consensus*, i.e.
majority vote

Issues with overlap graph assembly

- Next-generation sequencing produces 100's of millions (or even billions) of reads
- With one node per read this is computationally intractable for large genomes
- What if the nodes in our graph were not reads?

De Bruijn Graph (DBG) Assembly

DBG: de Bruijn Graph assembly

k-mer: substring of length k

S : G G C G A T T C A T C G

4-mer: A T T C

total
k-mers:

all 3-mers:

G G C

G C G

C G A

G A T

... T C G

read

A A A B B B A

k=3

① list all k-mers

↑
B
↑
A

x2

AAA

AAB

ABB

BBB

BBA

AA, AA

AA, AB

AB, BB

BB, BB

BB, BA

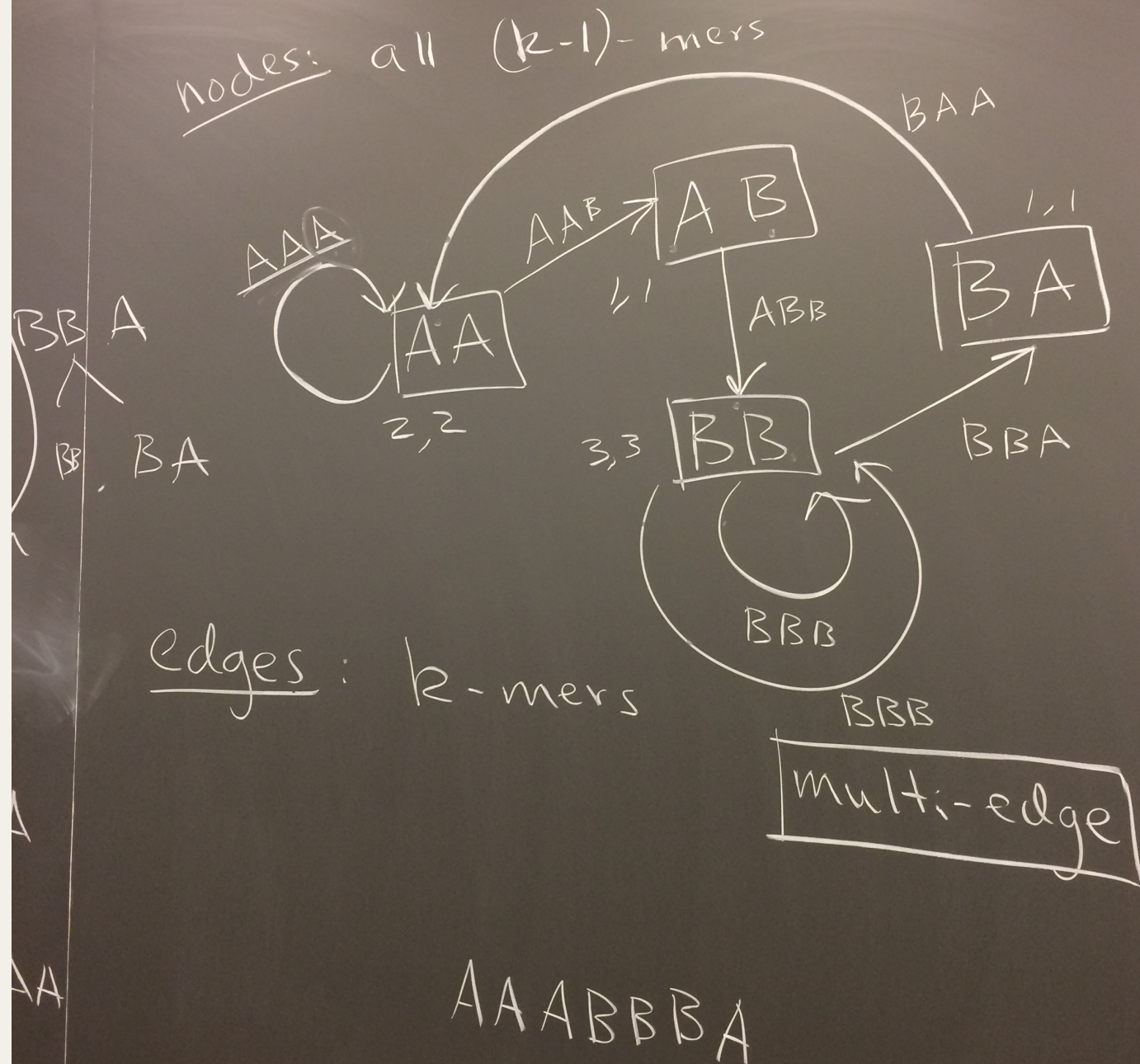
2^k

L/R (k-1)-mers

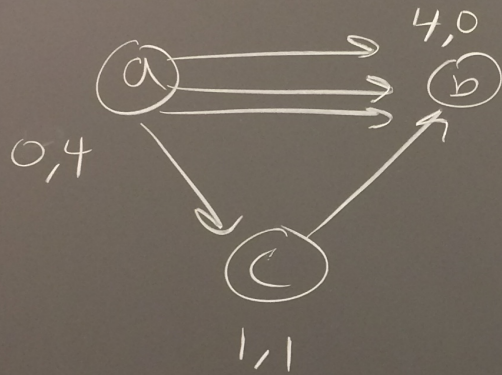
BAA

BA AA

ed



de Bruijn graph: directed multigraph.



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

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

indegree: # incoming edges

outdegree: # outgoing edges

node is balanced if: $\text{indegree} = \text{outdegree}$

semi-balanced: $|\text{indegree} - \text{outdegree}| = 1$

Goal: visit each edge exactly once

Eulerian path

k

Theorem:

There exists an Eulerian path iff all nodes are balanced except for 2 that are semi-balanced.

Z A B C D A B E F A B Y

$k=3$, what is the de Bruijn graph.