# COMP90014

Algorithms for Bioinformatics

Week 8A: Considerations & Graph Simplification

# Assembly Algorithms

# Recap

## Overlap Layout Consensus
Strategy of choice for long read assembly

| Overlap | Build overlap graph |
|---------|---------------------|
| Layout | Group parts of graph into contigs |
| Consensus | Squash reads to generate consensus contigs |

Overlap-Layout-Consensus (OLC)

Commins et al., 2009

Overlap          reads
       reads

Layout                    reads
            reads

Consensus

# Recap

## Overlap Layout Consensus

Strategy of choice for long read assembly

| | |
|---|---|
| **Overlap** | Build overlap graph |
| **Layout** | Group parts of graph into contigs |
| **Consensus** | Squash reads to generate consensus contigs |

Graph nodes = reads
Graph edges = overlaps

```
alignment
    overlap
```

## Overlap, Layout



simplification    traversal

| Size matt | not. Look | e. Jud | my soze, do |
| tters not. Lo | | me. Judge me by my size, d | |
| s not. Look at | | | size, do you? |
| | k at me. Judg | | |

Size matters not. Look at me. Judge me by my size, do you?

## Consensus

**reads:**
ATCGATGCTAGCTGA---
-----TGCTAGCTGATGA
-TCGA**A**G-TAG-TGATG-
AT**A**GATGCTAGCTGA-GA

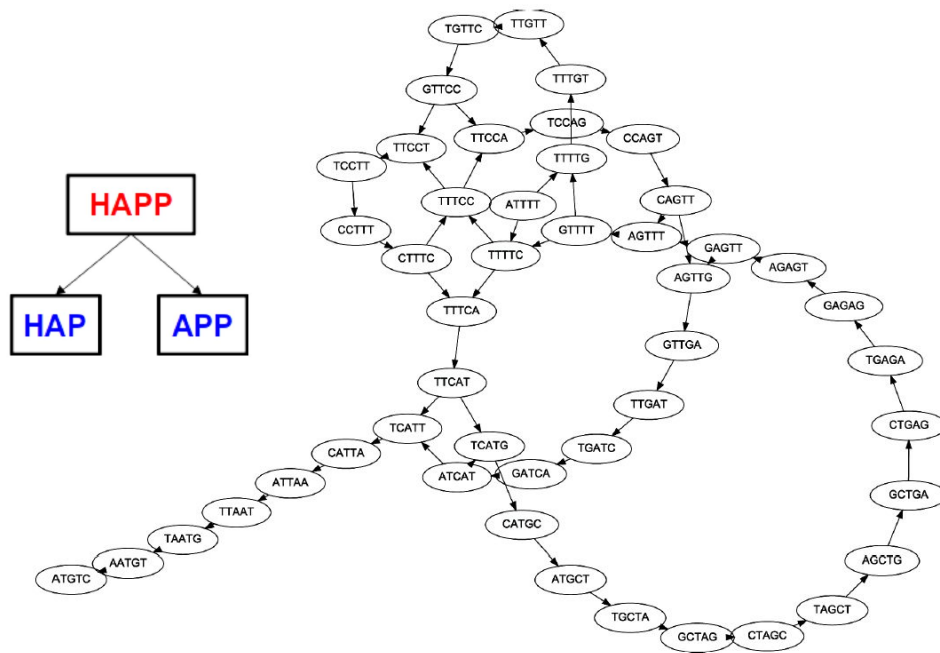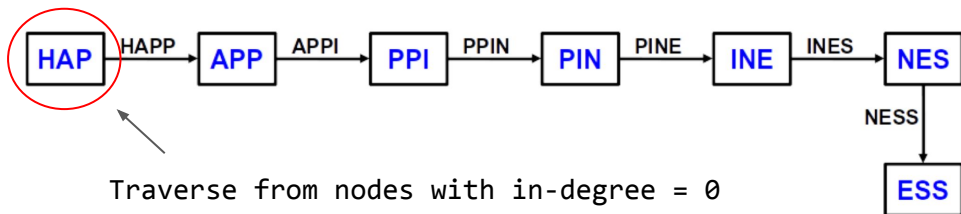**consensus:**
ATC GAT GC TAG CTGA TGA

# Recap

## De Bruijn Graphs

Strategy of choice for short read assembly

1. Break reads into kmers
2. Separate kmers into prefix - suffix
3. Add prefix / suffix and edge to graph

Graph nodes = prefixes / suffixes
Graph edges = kmers

# Recap

## De Bruijn Graphs

Strategy of choice for short read assembly

1. Break reads into kmers
2. Separate kmers into prefix - suffix
3. Add prefix / suffix and edge to graph

Graph nodes = prefixes / suffixes
Graph edges = kmers



Traverse from nodes with in-degree = 0

```
HAPP
 APPI
  PPIN
   PINE
    INES
     NESS
---------
HAPPINESS
```

The reads assemble to
HAPPINESS
Easy!

# Assembly Algorithms

Recap

OLC vs De Bruijn Assembly

Considerations

OLC Simplification

# OLC vs De Bruijn

|  | OLC | De Bruijn |
|---|---|---|
| Preprocessing | All-vs-all alignment (slow) | Kmer counting (fast) |
| Graph nodes | Reads | Kmer prefixes / suffixes |
| Graph edges | Overlaps | Kmers |
| Traversal | Hamiltonian (slow) | Eulerian (fast) |
| Read coherence | Yes | No |
| Repeat resolution | Better | Worse |
| Read errors | No problem | A problem |

Almost like OLC is well suited to long reads, and De Bruijn is suited to short reads…

# OLC vs De Bruijn

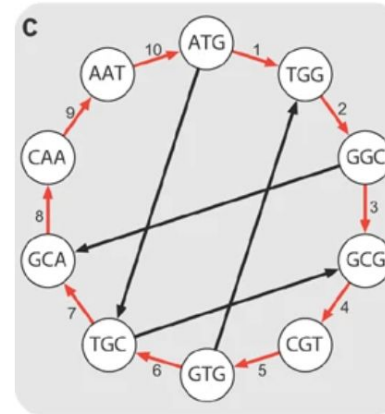Both OLC and De Bruijn approaches use graphs as intermediate.

OLC has some pros:

- Better tolerance of read errors
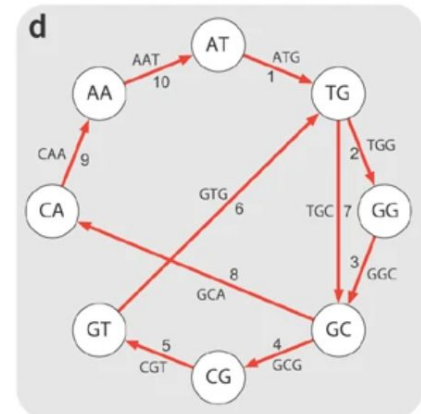- Better repeat resolution power

But enormous complexity hit!

- All vs all alignment (get overlaps)
- Hamiltonian vs Eulerian cycles

Why is De Bruijn much more scalable?



Hamiltonian cycle
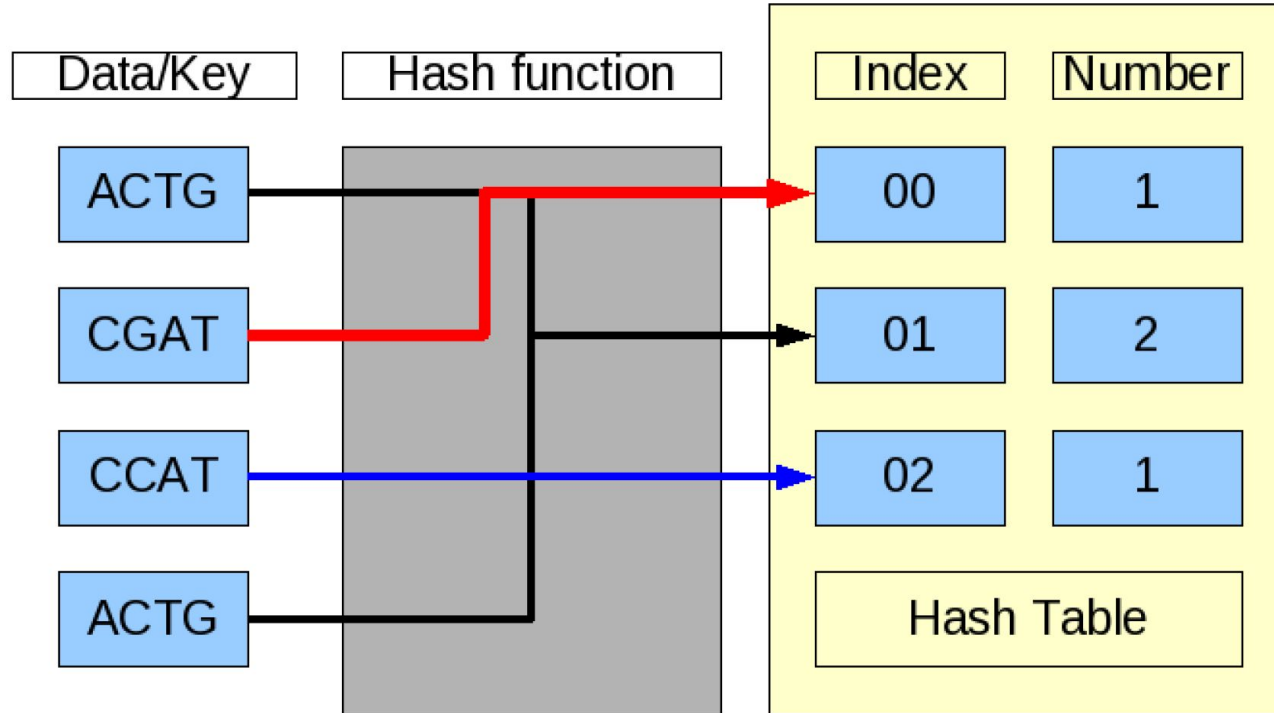Visit each vertex once
(harder to solve)

Eulerian cycle
Visit each edge once
(easier to solve)

Compeau, Pevzner, Tesler (2011)

# De Bruijn Efficiency

1. Select a value for *k*
2. **Hash** the reads (break them into *k*-mers)
3. Count the *k*-mers
4. Make the de Bruijn graph and annotate with *k*-mer frequency
5. Simplify the graph
6. Read off contigs from simplified graph

# De Bruijn Efficiency

# De Bruijn Efficiency

🍪 Make k-mers in $O(N_r)$ time ($N_r$ = number of reads)

🍪 Using de Bruijn graphs is fast

- exact matches: **no alignments to compute**!
- look up $k$-mers in index (**hash**) in linear time
- can build graph in $O(N_k)$ time ($N_k$ = number of unique $k$-mers)

🍪 Eulerian pathing is fast

- make contigs in $O(E)$ time ($E$ = number of edges)

# De Bruijn Efficiency

🍪 What de Bruijn graph advantages have we discovered?

- Can be built in $O(N)$ expected time
  $N$ = total number of reads
- With perfect data, graph uses $O(\min(N, G))$ space
  $G$ = genome length
- Note: when depth is high, $G \ll N$

🍪 Compares favourably with overlap graph:

- Fast overlap graph construction (suffix tree) is $O(N + a)$ time
  $a$ = number of alignments
- Space is $O(N + a)$.
- **Remember**: $a$ is $O(N^2)$

# Assembly Algorithms

Recap

OLC vs De Bruijn Assembly

<u>Considerations</u>

OLC Simplification

# Considerations

Last lecture - Toy examples.

In practice, graphs get very complex due to many considerations

# Considerations

Last lecture - Toy examples.

In practice, graphs get very complex due to many considerations

OLC specific
-   Graph has lots of edges
    Intrinsic to overlap graphs

De Bruijn specific
-   Effect of kmer size, use of kmers
-   Read errors
-   Ploidy

Graph has lots of edges

OLC                                                              reads
read                           reads


                    read                    reads


kmer                    Effect of kmer size
De Bruijn                    kmers      kmer           k                                              k
                                                        k


    kmers   use of kmers
        De Bruijn       reads                 k              kmers              kmers              kmer
                        kmer


        Read errors
                            reads                              De Bruijn


Ploidy

# Considerations

Last lecture - Toy examples.

In practice, graphs get very complex due to many considerations

OLC specific
- Graph has lots of edges
  Intrinsic to overlap graphs

De Bruijn specific
- Effect of kmer size, use of kmers
- Read errors
- Ploidy

Both
- Repeats
- Strandedness
- Will mention, then promptly forget these.

Repeats

reads                                                    read
                                                         OLC      De
Bruijn

Strandedness
            DNA                                   reads
                    reads                                  "
        "
reads

# OLC Specific

## Many edges

Intrinsic to overlap graphs.
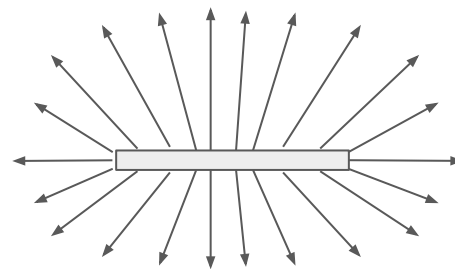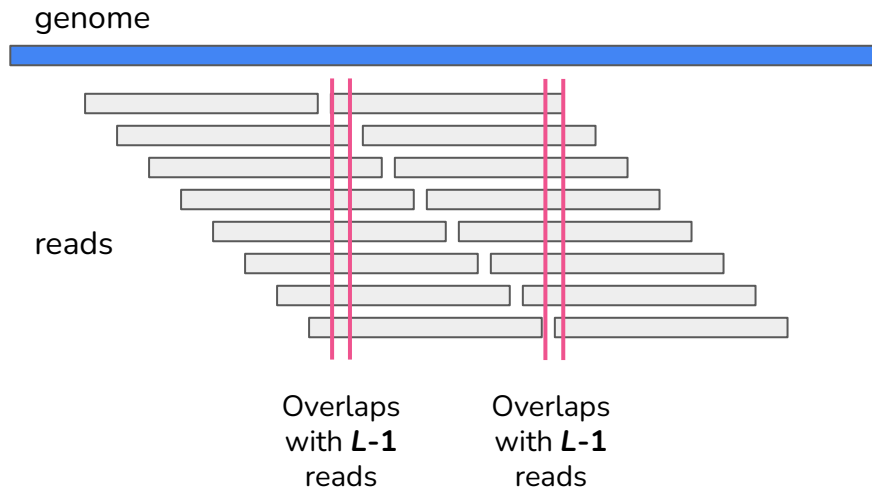
# OLC Specific

## Many edges

Intrinsic to overlap graphs.

- Sequencing depth = 30x
- Average 30 reads cover each position
- Each read overlaps with ~60 other reads
- Not including repetitive sequence
- This is a *big omission*



genome

reads

Overlaps with *L-1* reads

Overlaps with *L-1* reads

Sequencing depth = 30x　　　　　　　= 30　　　:
　　　　　　　　　　　　　　　　　　　30

Average 30 reads cover each position　　　　　　　30　reads　　　　:
　　　　　　　　　　　　　　　　　　　　　　　　　　30　　　　　reads

Each read overlaps with ~60 other reads　　　read　　　60　　　reads　　　:
　　　　　30　reads　　　　　　　　reads　　　　　read　　　　　　60　reads

Not including repetitive sequence　　　　　　　　:
　　　　　　　　　　　　　　　　　　　　　　　　read　　　60　　　　reads

This is a big omission　　　　　　　　:

# OLC Specific

## Many edges

Intrinsic to overlap graphs.

- Sequencing depth = 30x
- Average 30 reads cover each position
- Each read overlaps with ~60 other reads
- Not including repetitive sequence
- This is a **big omission**

Problematic when calculating hamiltonian path.

Can mitigate by enforcing min overlap len.

Some edges are redundant?

genome

reads

Overlaps
with **L-1**
reads

Overlaps
with **L-1**
reads

# De Bruijn Specific

## Kmers

Short reads already short!

    Reducing them further is a sin.. but no choice

    Lose repeat resolution power

    Lose read coherence

# De Bruijn Specific

## Kmers

Short reads already short!

   Reducing them further is a sin.. but no choice

   Lose repeat resolution power

   Lose read coherence

**Lower K:** Fewer nodes, more edges

**Higher K:** More nodes, fewer edges

Balance these two forces. See more Lecture 8B.

# De Bruijn Specific

## Ploidy

Heterozygosity causes many parallel splits and joins.

pat.

mat.

# De Bruijn Specific

## Ploidy

Heterozygosity causes many parallel splits and joins.

OLC: Reads from different strands overlap. Will be merged eventually.

De Bruijn: No notion of overlaps. Will need to correct these at some point...



pat.

mat.

# De Bruijn Specific

## Read Errors

Mainly affects De Bruijn Graphs

OLC: alignment allows mismatches, gaps
De Bruijn: use of kmers & exact matching

# De Bruijn Specific

Read Errors

Mainly affects De Bruijn Graphs

OLC: alignment allows mismatches, gaps
De Bruijn: use of kmers & exact matching

Middle of read: causes bubbles

End of read: cause spurs

Use of kmers: causes erroneous edges

Bubbles

Spurs
(tips / dead ends)

Erroneous
edges

——OLC (Overlap-Layout-Consensus)    De Bruijn   ——

OLC: alignment allows mismatches, gaps  OLC                        :
    OLC           reads                                        reads


De Bruijn: use of kmers & exact matching  De Bruijn      k-mers            :
De Bruijn            k-mers              k-mers                                    k-mer

Middle of read: causes bubbles                    "      "   :
                              De Bruijn                        "      "


End of read: cause spurs                  "    "   :
                              De Bruijn              "    "

Use of kmers: causes erroneous edges      k-mers              :
                              k-mers     De Bruijn

# De Bruijn Specific

Read Errors

Example

```
HAPPI INESS APLIN PINET
          *         *
```

Unique 3-mers (k = 3):

```
HAP APP PPI
 INE NES ESS
  APL PLI LIN
    PIN NET
```
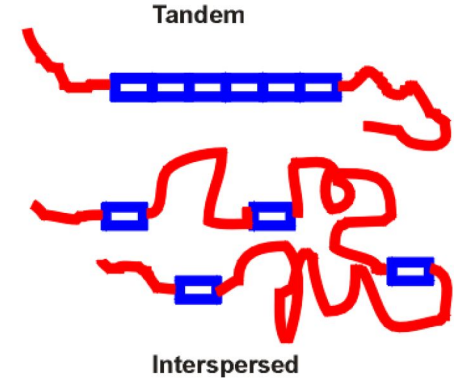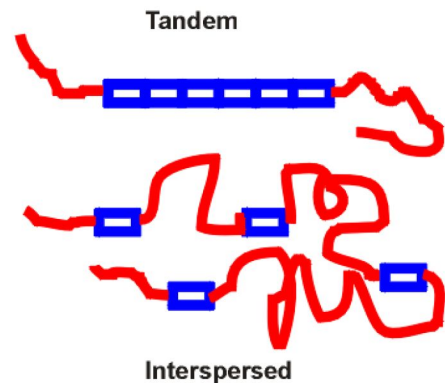
# Both - Repeats

Repeats

A segment of DNA that occurs more than once in the genome

Tandem repeats: Next to each other

Interspersed repeats: Have other sequence between repeats

Generally can't do much about repeats*

Repeats                    :

                                    DNA

Tandem repeats             :
                        DNA                                               DNA

Interspersed repeats            :
                        DNA                    DNA


*Generally can't do much about repeats                              **:

# Both - Repeats

## Repeats

A segment of DNA that occurs more than once in the genome

Tandem repeats: Next to each other

Interspersed repeats: Have other sequence between repeats

Generally can't do much about repeats*



| Repeat Class | Arrangement | Coverage (Hg) | Length (bp) |
|---|---|---|---|
| **Microsatellite** | Tandem | 3.00% | 2–100 |
| **SINE** | Interspersed | 15.00% | 100–300 |
| **Transposable elements** | Interspersed | 12.00% | 200–5k |
| **LINE** | Interspersed | 21.00% | 500–8k |
| **rDNA** | Tandem | 0.01% | 2k–43k |
| **Segmental Duplications** | Tandem or Interspersed | 0.20% | 1k–100k |

# Both - Repeats

## Tandem repeats

Tandem repeat

Short reads

# Both - Repeats

## Tandem repeats

Tandem repeat

Short reads

1x? 10x? 100x?

De Bruijn Graph

# Both - Repeats

## Tandem repeats



Tandem repeat

Short reads

Long reads

1x? 10x? 100x?

De Bruijn Graph

# Both - Repeats

## Tandem repeats



Tandem repeat

Short reads

Long reads

1x? 10x? 100x?

De Bruijn Graph

Overlap Graph

# Both - Repeats

## Interspersed repeats

Copy 1

Copy 2

Short reads

# Both - Repeats

Interspersed repeats

# Both - Repeats

## Interspersed repeats



Copy 1  Copy 2

Short reads

Long reads

De Bruijn Graph

# Both - Repeats

# Both - Strandedness
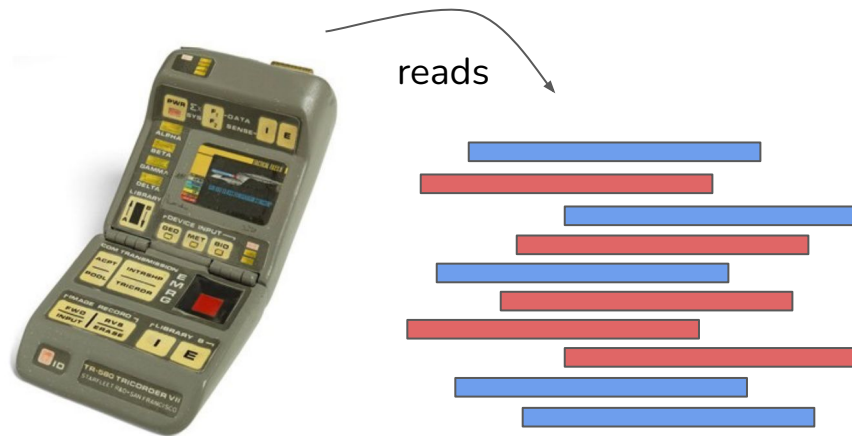
DNA and English sentences are just strings of characters

# Both - Strandedness

DNA and English sentences are just strings of characters

5' ATCGATCGTAAATGGC 3'
3' TAGCTAGCATTTACCG 5'

DNA is double stranded
🍪 minus strand is the *reverse complement* of the plus strand

OLC:

OLC

reads        reads

DNA        reads

reads    strandedness        reads

strandedness

De Bruijn    :

De Bruijn        k-mers

k-mer        k-mer

De Bruijn        reads

k-mers        strandedness

# Both - Strandedness

DNA *and English sentences are just strings of characters*

5' ATCGATCGTAAATGGC 3'
3' TAGCTAGCATTTACCG 5'

DNA is double stranded
- 🍪 minus strand is the **reverse complement** of the plus strand
- 🍪 (most) sequencing is **NOT STRAND SPECIFIC**
- 🍪 how do we deal with this in the graph?

reads

# Both - Strandedness - OLC

Overlap graphs (seen in OLC) use additional data on the edges.

# Both - Strandedness - OLC

Overlap graphs (seen in OLC) use additional data on the edges.

Reads:

```
(f)  TTAC
(f)   TACGA
(r)    ACGAA -> TTCGT
(f)      GAAT
(f)       AATATA
```
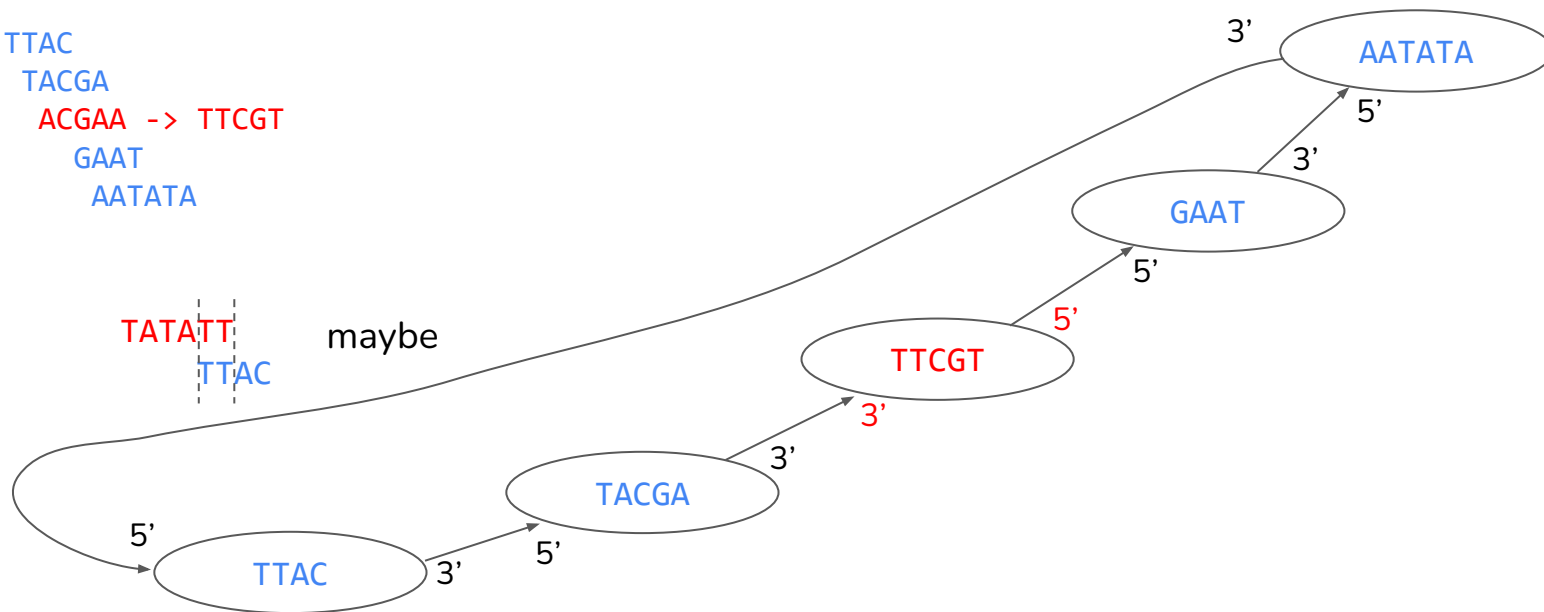
# Both - Strandedness - OLC

Overlap graphs (seen in OLC) use additional data on the edges.

Reps:

(f) TTAC
(f)  TACGA
(r)   ACGAA -> TTCGT
(f)     GAAT
(f)      AATATA

# Both - Strandedness - OLC

Overlap graphs (seen in OLC) use additional data on the edges.

Reads:

(f)  TTAC
(f)   TACGA
(r)    ACGAA -> TTCGT
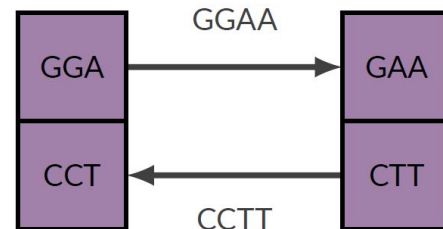(f)      GAAT
(f)       AATATA

# Both - Strandedness - OLC

Overlap graphs (seen in OLC) use additional data on the edges.

Reads:

(f) TTAC
(f)  TACGA
(r)   ACGAA -> TTCGT
(f)     GAAT
(f)      AATATA

TATATT     maybe
 TTAC

# Both - Strandedness - De Bruijn

De Bruijn Graphs directly add the reverse-complemented kmers to the graph

# Both - Strandedness - De Bruijn

De Bruijn Graphs directly add the reverse-complemented kmers to the graph

1. Take the first *k*-mer
2. Split into prefix/suffix
3. Add them as nodes
4. Add the reverse complement to nodes
5. Add *k*-mers (forwards and reverse complement) as edges
6. Repeat for the rest of the *k*-mers



```
GGAA   GAAC   AACT   TTCC   TTCG   GTTC   AGTT
^^^^
```

# Both - Strandedness - De Bruijn

De Bruijn Graphs directly add the reverse-complemented kmers to the graph

GGAA   GAAC   AACT   TTCC   TTCG   GTTC   AGTT



contigs

GGA   GAAC   CGAA

Traversing the graph: *move in the forwards direction only*

Contigs: GGAA   GAACT   CGAA

# Summary

Our graphs get very complex! Need to simplify.

# Summary

Our graphs get very complex! Need to simplify.

Overlap Graphs (seen in OLC)

Unnecessary edges
- Intrinsic to overlap graphs

Graph has spurs
- Due to sequencing errors chimeric sequences (repeats)

# Summary

Our graphs get very complex! Need to simplify.

**Overlap Graphs (seen in OLC)**

Unnecessary edges
- Intrinsic to overlap graphs

Graph has spurs
- Due to sequencing errors chimeric sequences (repeats)

**De Bruijn Graphs**

Graph has bubbles
- Read errors (middle of read), heterozygosity

Graph has spurs
- Read errors (end of read)

Graph has erroneous edges
- Use of kmers

# Assembly Algorithms

Recap

OLC vs De Bruijn Assembly

Considerations

OLC Simplification

# Overlap Graph Simplification

Want a hamiltonian path

   ...but overlap graphs have many edges & dead ends

# Overlap Graph Simplification

Want a hamiltonian path

...but overlap graphs have many edges & dead ends
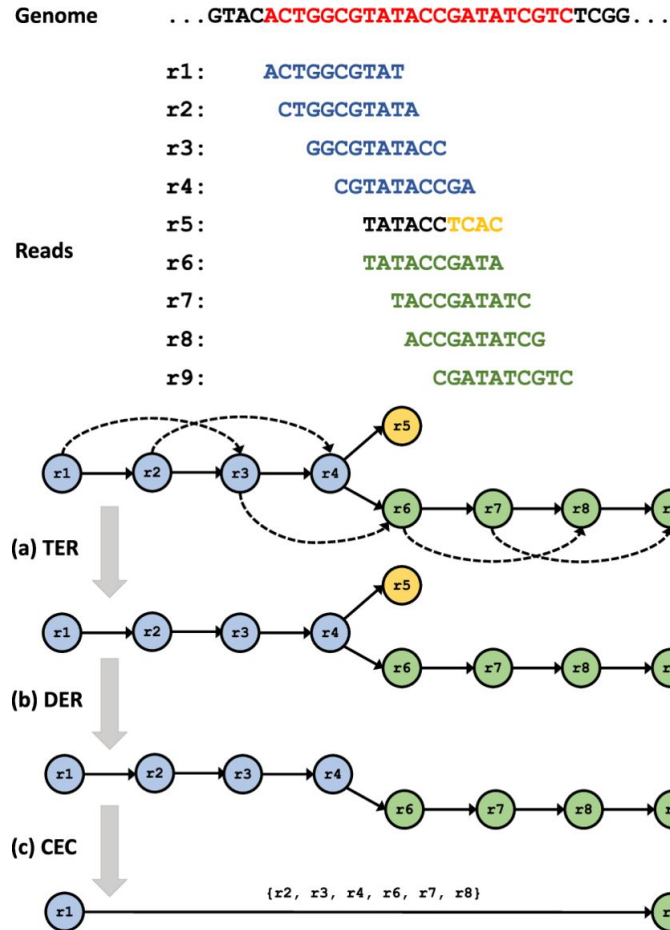
Transitive Edge Reduction (TER)

- Remove unnecessary edges



Paul et al (2019)

# Overlap Graph Simplification

Want a hamiltonian path

...but overlap graphs have many edges & dead ends

Transitive Edge Reduction (TER)

- Remove unnecessary edges

Dead-End Removal (DER)

- Remove short spurs / dead ends



Paul et al (2019)

# Overlap Graph Simplification

Want a hamiltonian path

...but overlap graphs have many edges & dead ends

Transitive Edge Reduction (TER)

- Remove unnecessary edges

Dead-End Removal (DER)

- Remove short spurs / dead ends

Composite Edge Contraction(CEC)

- Merges nodes in manner which does not lose information
- Quite complex. Not covering this.



Paul et al (2019)

Want a hamiltonian path

read

...but overlap graphs have many edges & dead ends

Transitive Edge Reduction (TER)

Remove unnecessary edges
TER

Dead-End Removal (DER)

Composite Edge Contraction (CEC)

# Transitive Edge Reduction

Overlap graphs inherently messy.

Each read overlaps with many other reads.

to_every_thing_turn_turn_turn_there_is_a_season

Read length = 7
Min overlap = 4

# Transitive Edge Reduction

Luckily, many of the overlaps are redundant and can be removed.

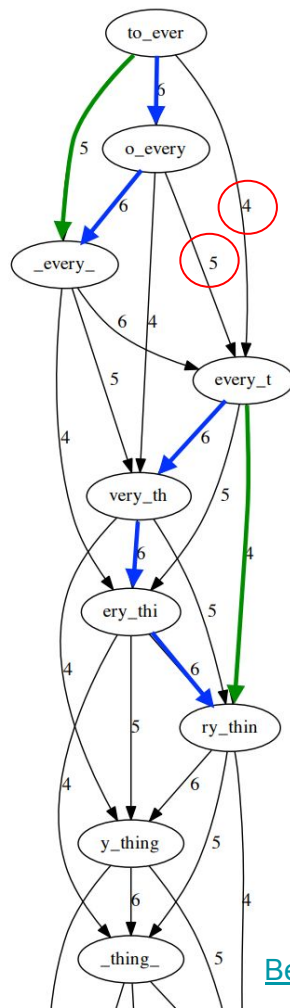Call these edges *transitive*

# Transitive Edge Reduction

Luckily, many of the overlaps are redundant and can be removed.

Call these edges *transitive*

Transitive edges can be inferred from other edges.

In the figure, **green** edges can be inferred from **blue** edges.

Know which path to walk when removing these due to the **overlap lengths**

# Transitive Edge Reduction

Luckily, many of the overlaps are redundant and can be removed.

Call these edges *transitive*

Transitive edges can be inferred from other edges.

In the figure, **green** edges can be inferred from **blue** edges.

Know which path to walk when removing these due to the **overlap lengths**



**Algorithm 1** Algorithm for Transitive Edge Reduction

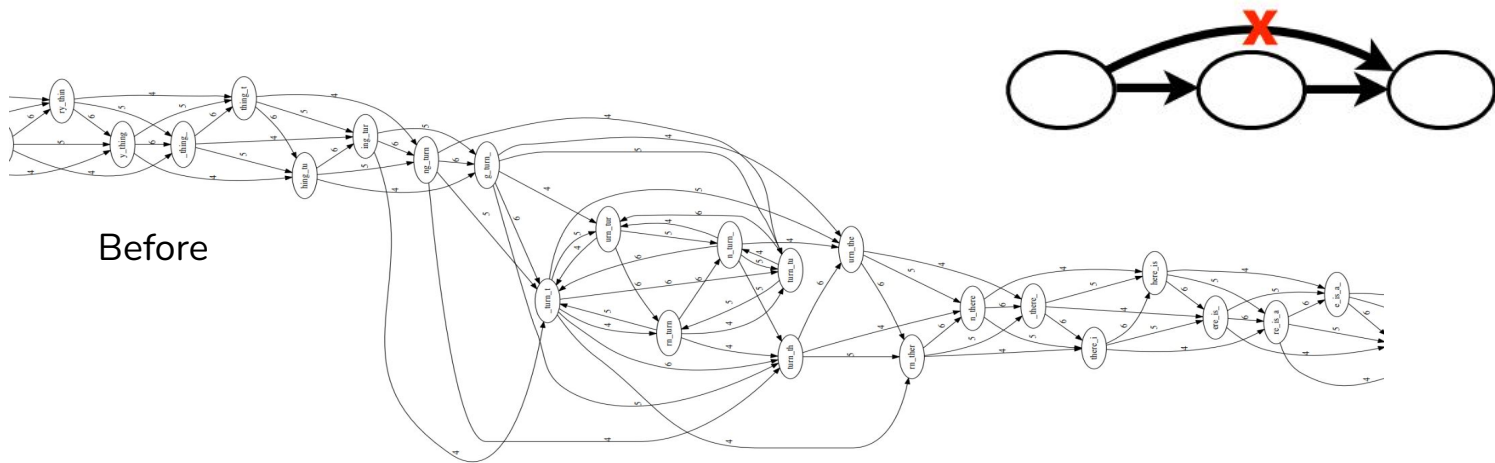**Input:** Graph (V, E)
**Output:** Reduced Graph (V, E')

1: **for** $v \in V$ **do**
2:  mark[v] ← vacant // Initially mark a vertex $v_i$ as vacant.
3:  **for** $v \rightarrow w \in E$ **do**
4:   reduce[$v \rightarrow w$] ← false // Mark an edge $w_i$ as not reduced.
5:  **end for**
6: **end for**
7: **for** $v \in V$ **do**
8:  // Mark a vertex $v_i$ reachable from $v_j$ as inplay
9:  **for** $v \rightarrow w \in E$ **do**
10:   mark[w] ← inplay
11:  **end for**
12:  longest ← $max_w len(v \rightarrow w)$
13:  **for** $v \rightarrow w \in E$ in order of length **do**
14:   // Traverse an edge $w_i$ marked inplay, indicating it is adjacent to $v$.
15:   **if** mark[w] = inplay **then**
16:    // Stop if an edge is too long to eliminate edges out of v.
17:    **for** $w \rightarrow x \in E$ in order of length **and** $len(w \rightarrow x) + len(v \rightarrow w) \leq longest$ **do**
18:     **if** mark[x] = inplay **then**
19:      mark[x] ← eliminated
20:     **end if**
21:    **end for**
22:   **end if**
23:  **end for**
24:  // Conclude the processing of $v$ by examining each adjacent vertex.
25:  **for** $v \rightarrow w \in E$ **do**
26:   **if** mark[w] = eliminated **then**
27:    reduce[$v \rightarrow w$] ← true // Mark as needing reduction.
28:   **end if**
29:   mark[w] ← vacant // Restore each vertex mark to vacant.
30:  **end for**
31: **end for**

Paul et al (2019)

Ben Langmead

# Transitive Edge Reduction

Process: (1) Remove edges which skip one node
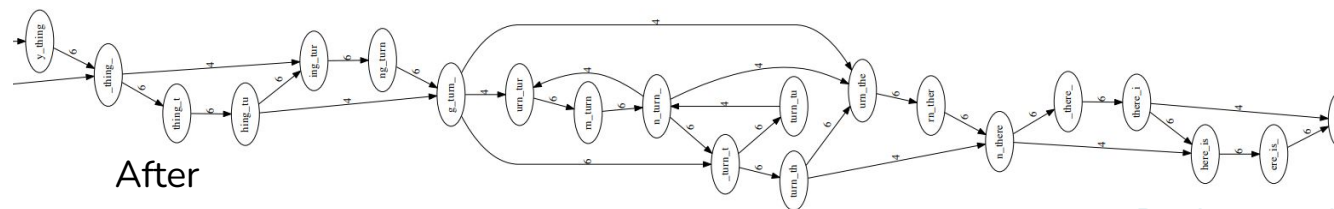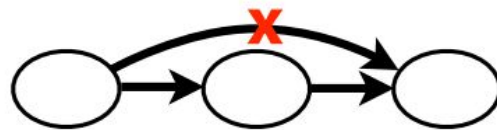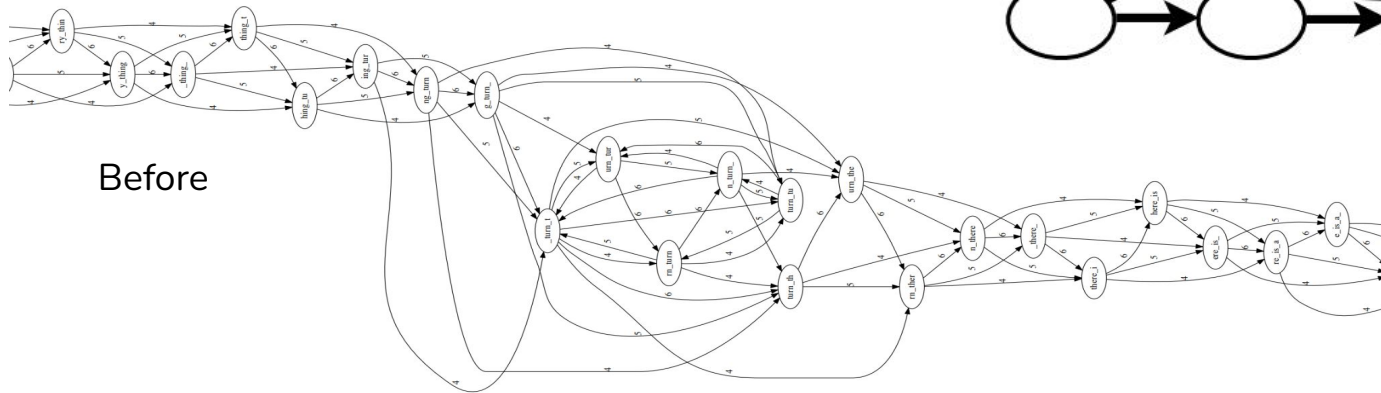


Before

Transitive Edge Reduction (TER)

Remove unnecessary edges
    TER                                          reads A   B    C      A      B B        C
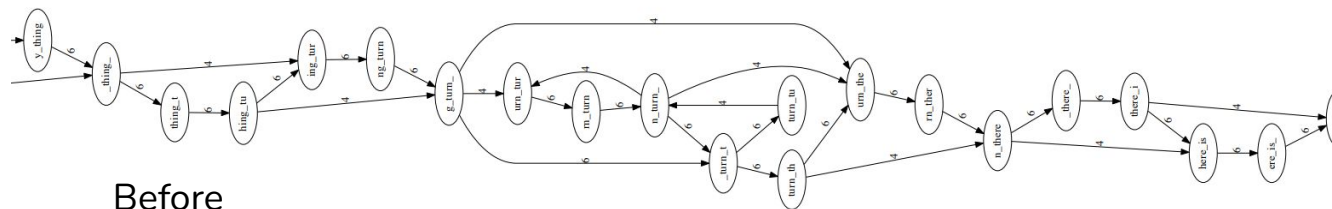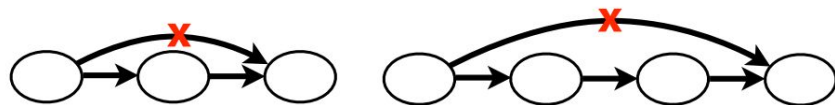        A              C            A   C                     A   B    C

# Transitive Edge Reduction

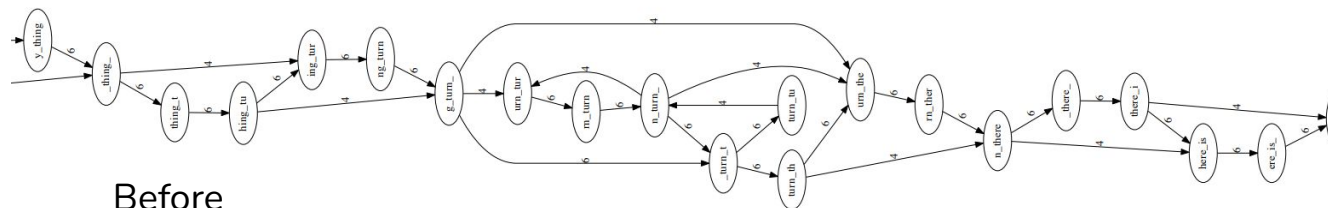Process: (1) Remove edges which skip one node



Before

After

# Transitive Edge Reduction

Process: (2) Remove edges which skip one or two nodes
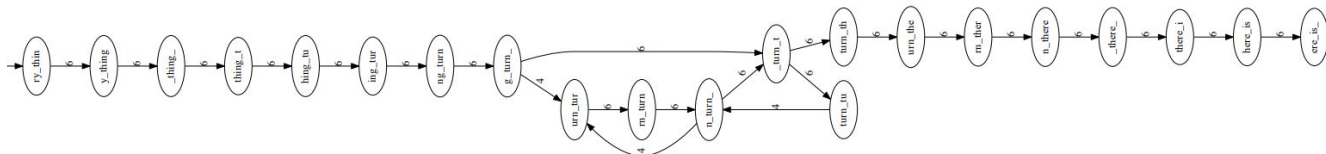


Before

# Transitive Edge Reduction

Process: (2) Remove edges which skip one or two nodes
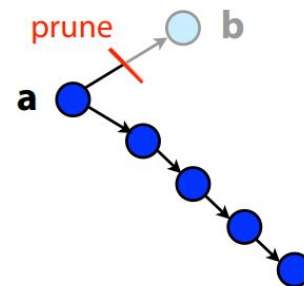


Before
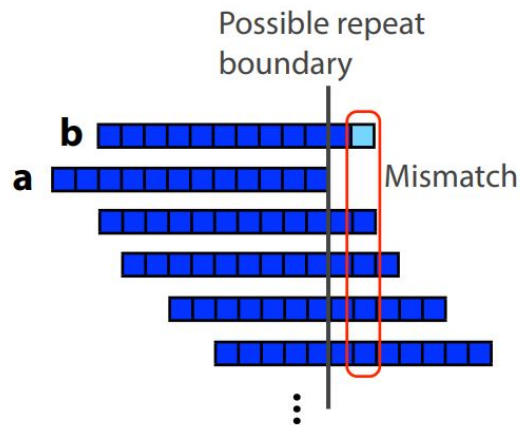
After

# Dead-End Removal (DER)

Remove short spurs / dead ends

Caused by sequencing errors

Caused by overlapping of chimeric sequences (repeats)

Simple to remove

- Identify, then prune
- Short length edges
- Low coverage (depth)



Ben Langmead

# Thank you!

**Today:** Considerations & Graph Simplification

**Next time:** Assembly in Practice