

Introduction to genome assembly

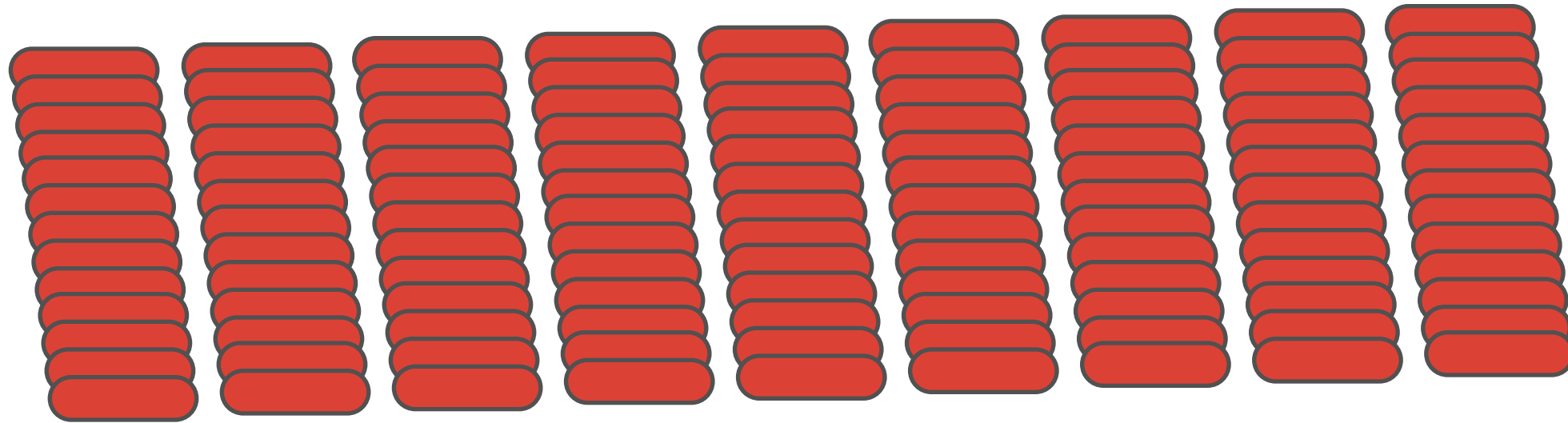
CMSC423

Many slides courtesy of Ben Langmead

SEC-GEN SEQUENCING

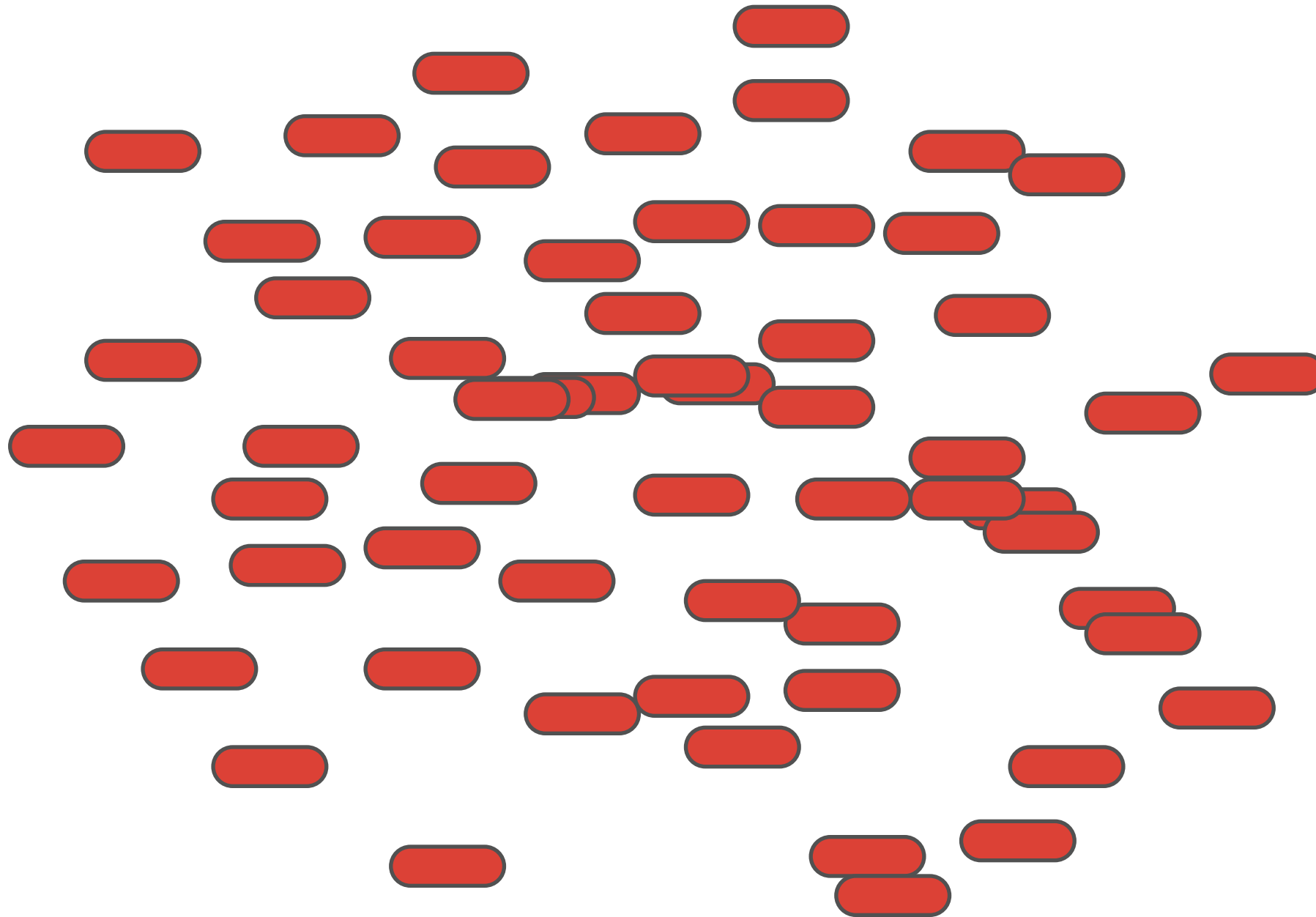


SEC-GEN SEQUENCING



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

SEC-GEN SEQUENCING



SECOND-GENERATION 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



```
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGGCGCTGNNNNNNNNNNCNCNNNN
+
BBBB>A?B@;@BBBBBAA=BA=A%????????????????????
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGCAGCAGGNNNNNNNNNGNNGNNNN
+
B9B@B<;BAA<@AB9=1>%????????????????????
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNCNCNNNCNCNNNN
+
A=B7&7:>B@:A>?9:<;:>?4?%????????????????????
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAGAAGCACAGCCAANNNNNANTNNCTNNNN
+
BBCCCCCBBBCB7CBC=7>+<=>=BCBCB%????????????????
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCCTCNNNNGNTNAAGNNNN
+
BCC?+<B=?BB5=ABA?B6BBBB4BB?B%????????????????
@HWI-EAS146:5:1:2:947#0/1
CCCAGGAGAAAGCCATGTTCAAGTTCGAGCGCANNANANCGTGANNNN
+
BBB9@?7A7>AAB@>?B=?@.>8?B?%????????????????
@HWI-EAS146:5:1:2:563#0/1
CCAGCCCCCTCCCCATCTCCCACCCTGTACCTNANCCCCTGANNNN
+
BBABAABB;AAABA77@5AAA:??>%????????????????
@HWI-EAS146:5:1:2:1631#0/1
TGGAACGCAGCCTACACTCTTCCAGGCCTCCTNCCTCCGTNNNN
+
BBB@6@BBBBBBBBB@BBBABAABBB?;9BB@BA5&<B:%?????
@HWI-EAS146:5:1:2:1420#0/1
CTCAAACCTCTGACCTTTGGTGATCCACCCGCCTNGGCCTTCNNNN
+
BBBB:BBBBBABAAA?:(=A8@>AAA?AB?=A%????????????
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCGAGGCTCCGCGGCGCTGNNNNNNNNNNCNCNNNN
+
BBBB>A?B@;@BBBBBAA=BA=A%????????????????????
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGCAGCAGGNNNNNNNNNGNNGNNNN
+
B9B@B<;BAA<@AB9=1>%????????????????????
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNCNCNNNCNCNNNN
+
A=B7&7:>B@:A>?9:<;:>?4?%????????????????????
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCAGAAGCACAGCCAANNNNNANTNNCTNNNN
+
BBCCCCCBBBCB7CBC=7>+<=>=BCBCB%????????????????
```

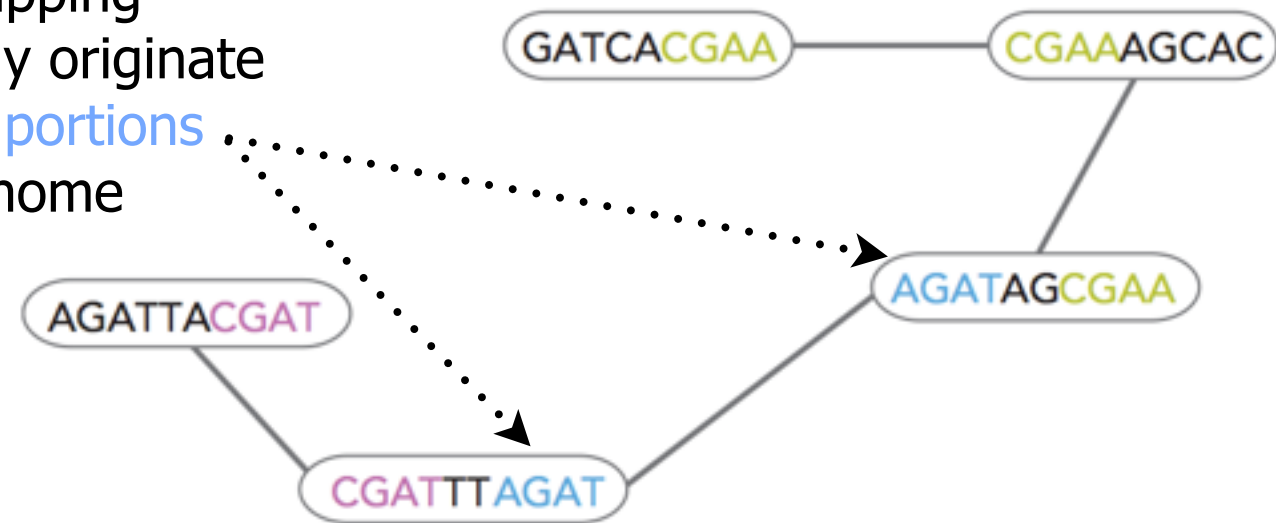
From reads to evidence

I. *de novo*

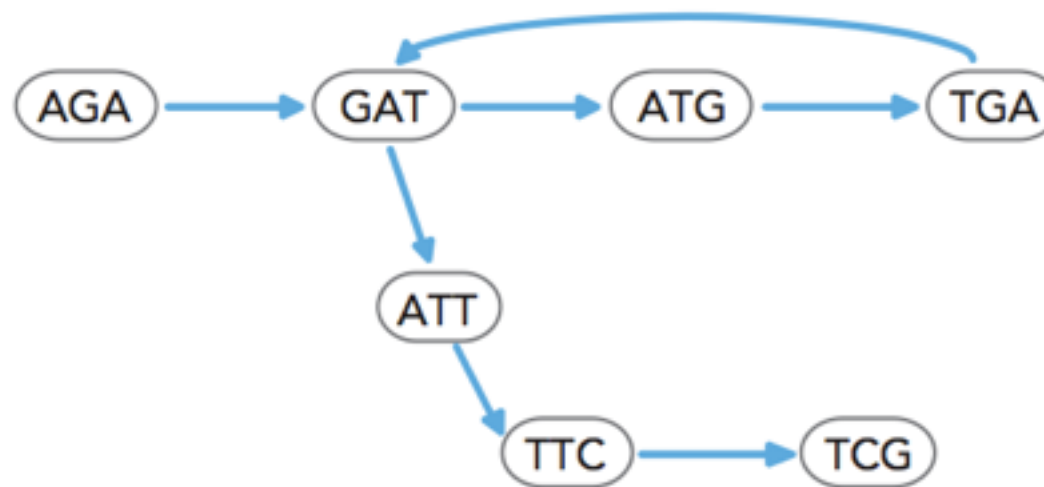
Assume nothing! - let reads tell us everything

Reads with overlapping sequence probably originate from **overlapping portions** of the subject genome

Encode overlap relationships as a graph



Source: De Novo Assembly Using Illumina Reads. Illumina. 2010



The full genome sequence is a "tour" of the graph

Source: De Novo Assembly Using Illumina Reads. Illumina. 2010

http://www.illumina.com/Documents/products/technotes/technote_denovo_assembly.pdf

```
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCAGGCTCCGCGGCGCTGNNNNNNNNNNNNNNNNNN
+
BBBB>A?B@;@BBBBBAA=BA=A%*****
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGCAGCAGNNNNNNNNNNNGNNNNNN
+
B9B@B<;BAA<@AB9=1>%*****
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNNNNNNNNNNNNN
+
A=B7&7:>B@:A>79:<;:>74?%*****
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCCAGAAGCACAGCCAANNNNNANTNNCTNNNN
+
BBCCCCCBB7CBC=7>+<=>BCBCB%*****
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCCTCNNNNNGNTNAAGNNNN
+
BCC?+<B=?BB5=ABA?B6BBBB4BB?B%*****
@HWI-EAS146:5:1:2:947#0/1
CCCAGGAGAAAGCCATGTTTCAGTTCGAGCGCANNANANCGTGANNNN
+
BBB9@?7A7>AAB@>?B=?@.>B?B?%*****
@HWI-EAS146:5:1:2:563#0/1
CCAGCCCCTCCCCATCTCCACCCTGTACCTNANCCCCTGANNNN
+
BBA8AABB;AAABA77@5AAA:??>%*****
@HWI-EAS146:5:1:2:1631#0/1
TGGGAACGCAGCCTACACTCTTCCAGGCCTCCTNCCTCCGTNNNN
+
BBB@6@BBBBBBBB@BBBABAABBB?;9BB@BA5&<B:%*****
@HWI-EAS146:5:1:2:1420#0/1
CTCAAACCTCTGACCTTTGGTGATCCACCGCCTNGGCCTTCNNNN
+
BBBB:BBBBBABAAA?: (=A8@>AAA?AB?=A%*****
@HWI-EAS146:5:1:1:961#0/1
TCCGAGGCCAACCAGGCTCCGCGGCGCTGNNNNNNNNNNNNNNNNNN
+
BBBB>A?B@;@BBBBBAA=BA=A%*****
@HWI-EAS146:5:1:1:1595#0/1
TCAGGAAGCAGGAAGAGCTGGTGCAGCAGNNNNNNNNNNNGNNNNNN
+
B9B@B<;BAA<@AB9=1>%*****
@HWI-EAS146:5:1:1:1048#0/1
CTGGACTGCATCCTACCACCAACTCGTCCAANNNNNNNNNNNNNNNN
+
A=B7&7:>B@:A>79:<;:>74?%*****
@HWI-EAS146:5:1:1:1607#0/1
CTCCTCTCAAGGTCCCCAGAAGCACAGCCAANNNNNANTNNCTNNNN
+
BBCCCCCBB7CBC=7>+<=>BCBCB%*****
@HWI-EAS146:5:1:1:1719#0/1
CACGATCTGGGTTTATTGTAACCTCCGCCTCNNNNNGNTNAAGNNNN
```

What we'll cover

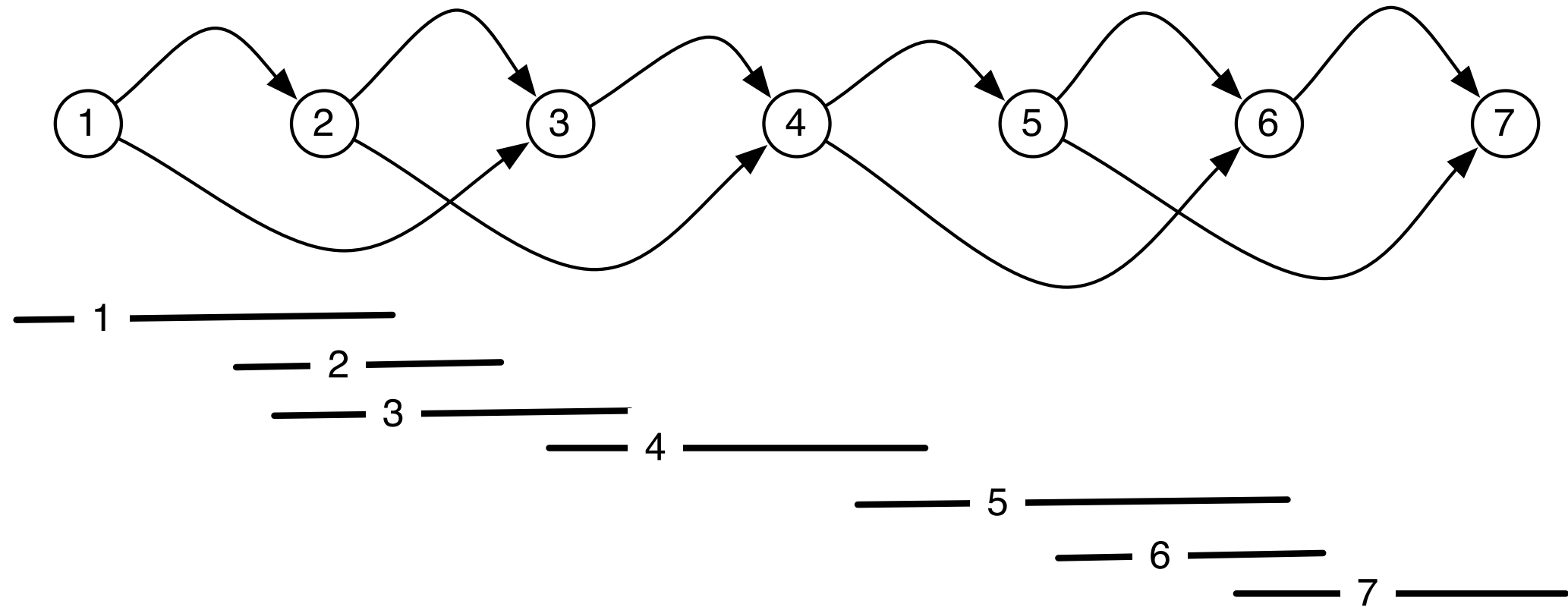
- Genome assembly as graph problems
 - Two representations:
 - Overlap graph
 - How much sequencing required for assembly
 - DeBruijn graph
- How to get assemblies from solutions to graph problems

Overlap Graph

Overlap graph:

Nodes = reads

Edges = overlaps



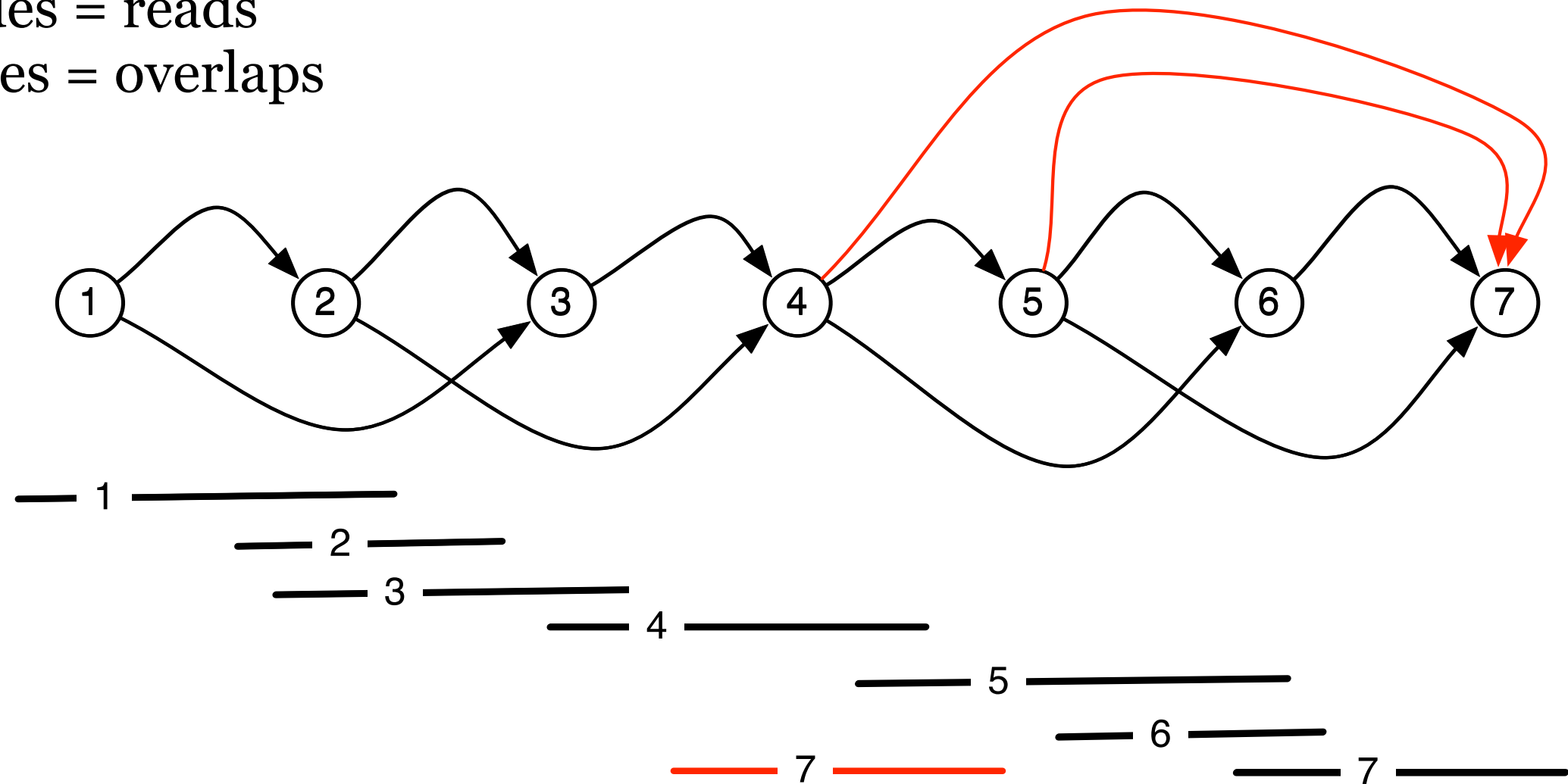
Given overlap graph, how can we find a good candidate assembly?

Overlap Graph

Overlap graph:

Nodes = reads

Edges = overlaps



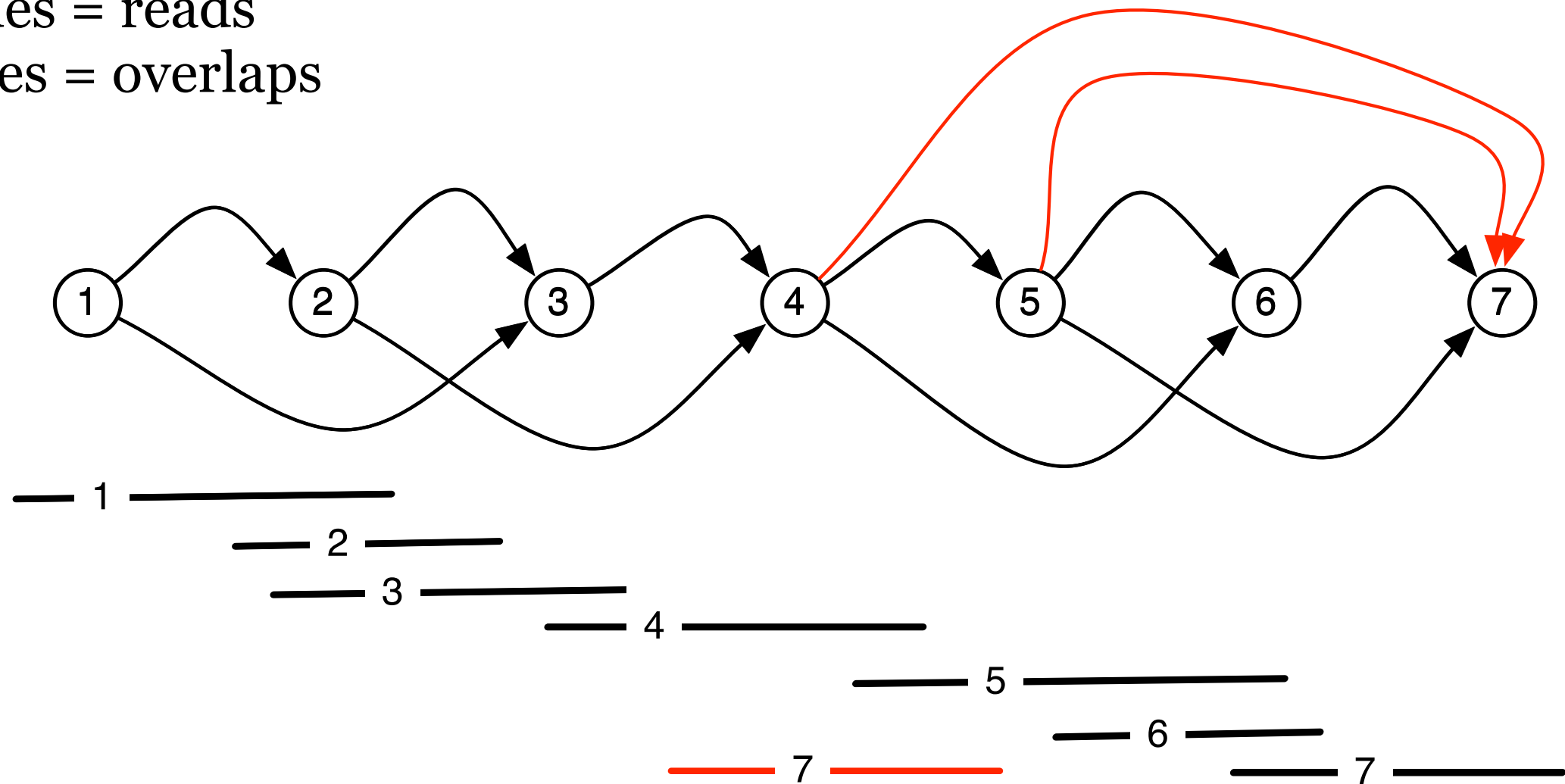
Given overlap graph, how can we find a good candidate assembly?

Overlap Graph

Overlap graph:

Nodes = reads

Edges = overlaps

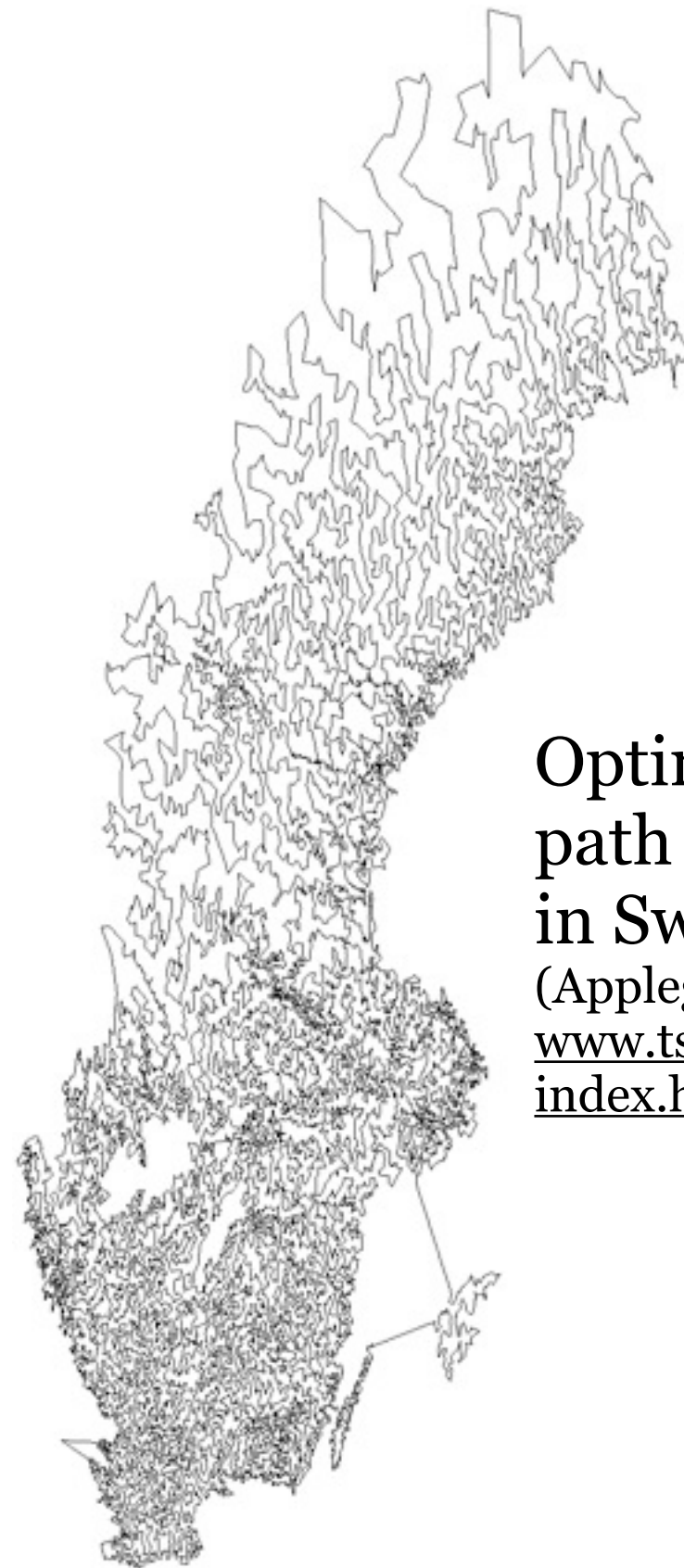


Given overlap graph, how can we find a good candidate assembly?

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.

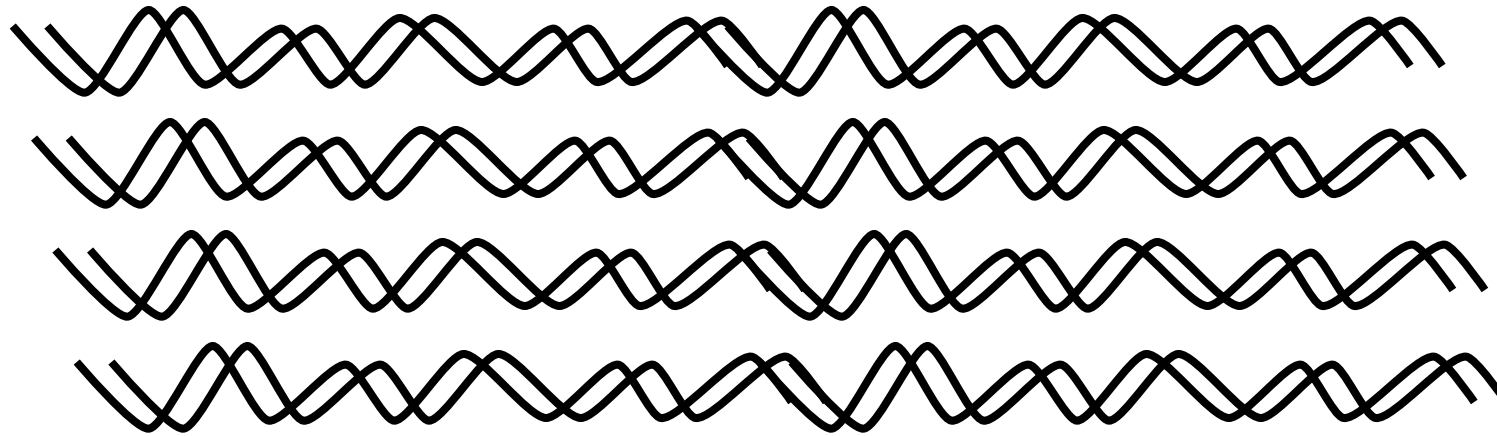


Optimal Hamiltonian path of 24,978 cities in Sweden

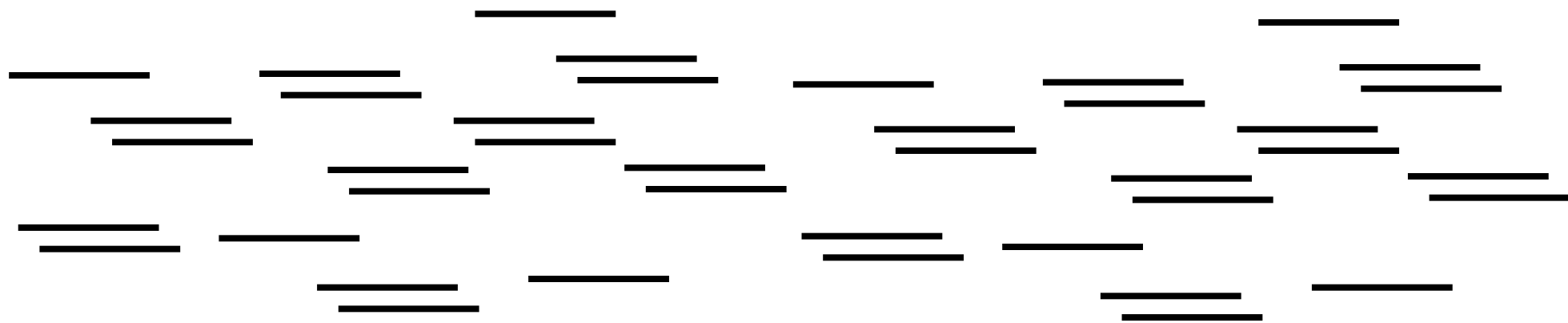
(Applegate et al, 2004, www.tsp.gatech.edu/sweden/index.html).

Shotgun Sequencing

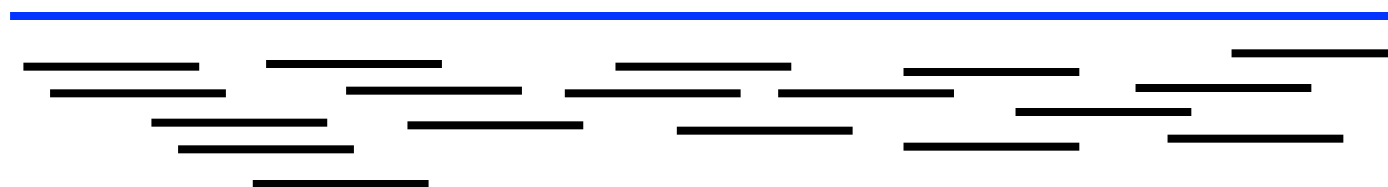
Many copies
of the DNA



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

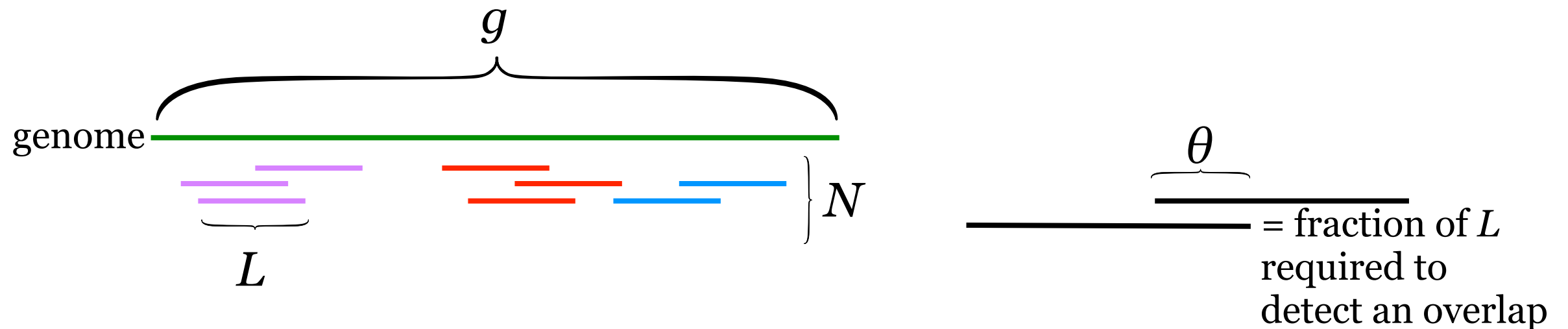


Assemble into original genome:



Lander-Waterman Statistics

How many reads do 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, θ .

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 λ of success

Expected # of successes = λx

Poisson approximation to binomial distribution:

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

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

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

Expected # of Islands, 2

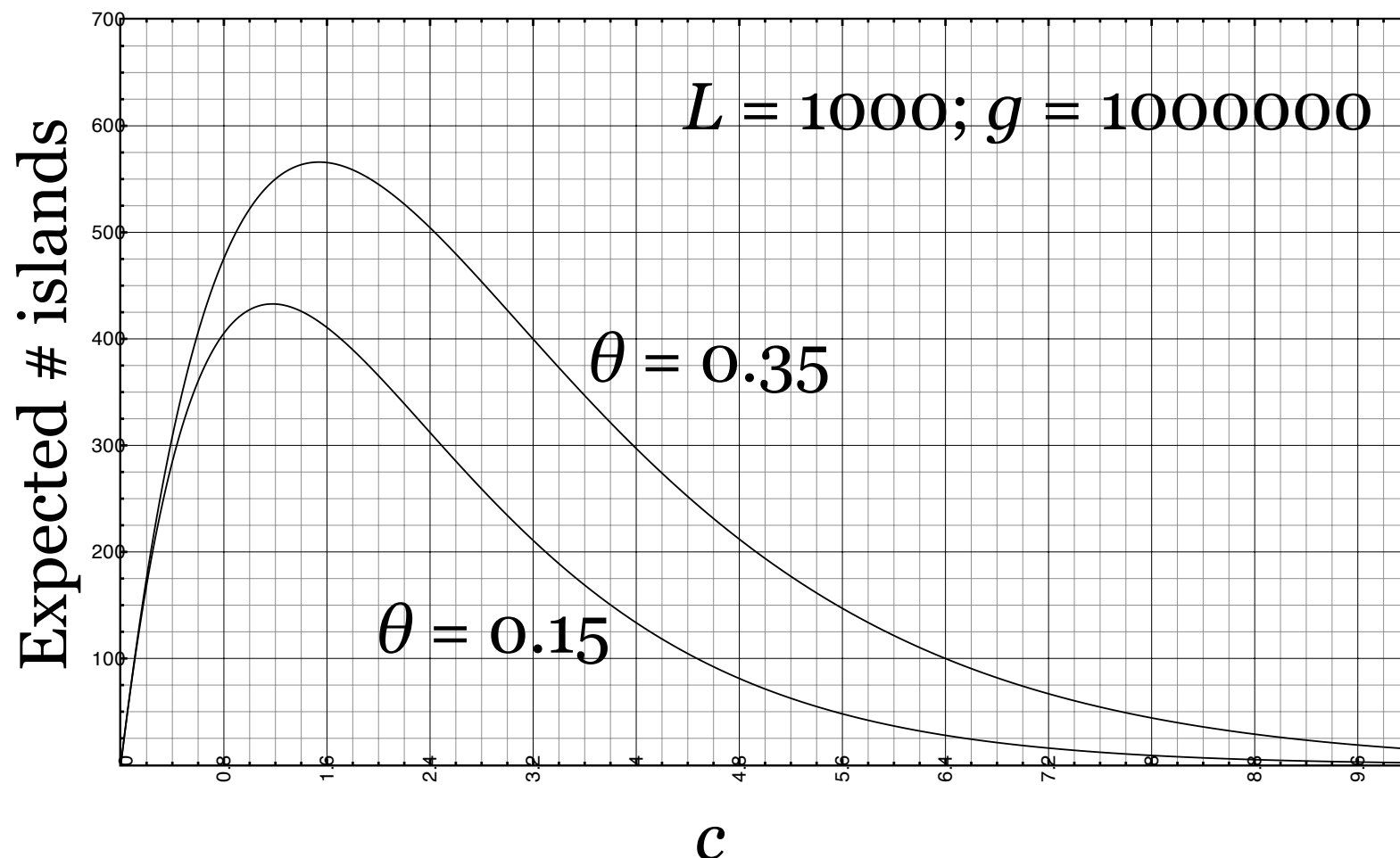
Rewrite to depend more directly on the things we can control: c and θ

$$\text{Expected \# of islands} = N e^{-(1-\theta) L N / g}$$

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

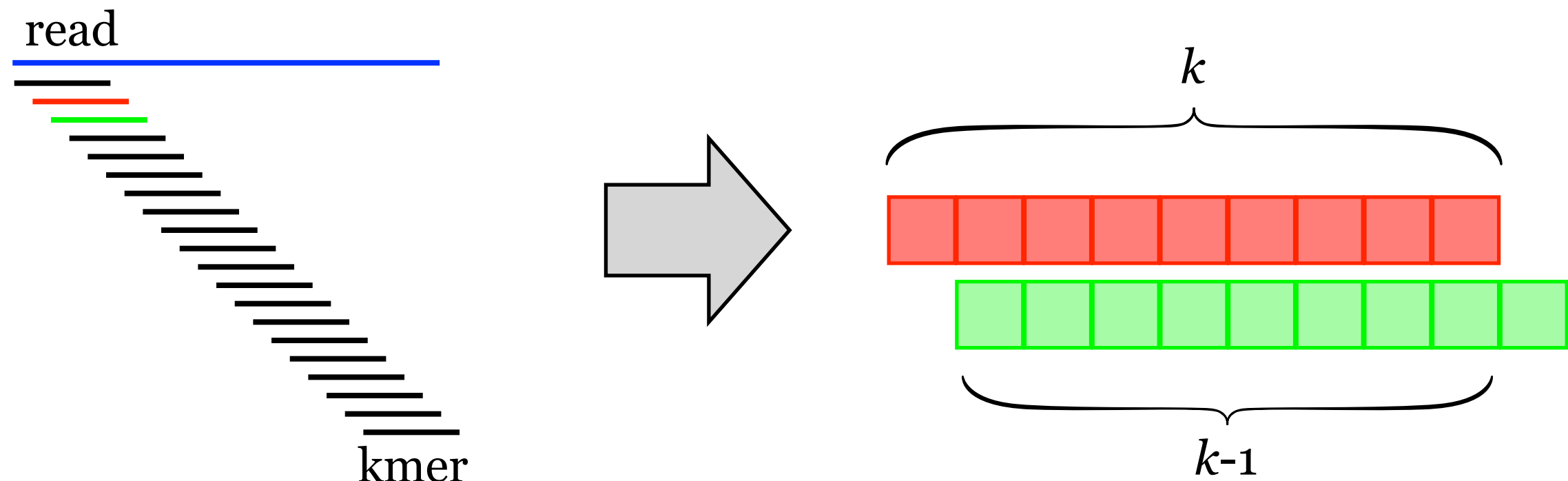
$$= \frac{L/g}{L/g} N e^{-(1-\theta) c}$$

$$= \frac{g}{L} c e^{-(1-\theta) c}$$



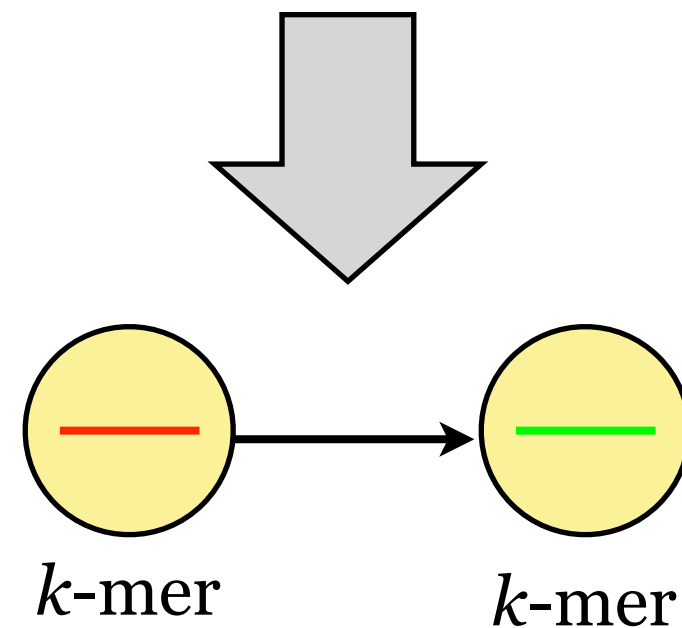
Assembly via Eulerian Path

de Bruijn graph

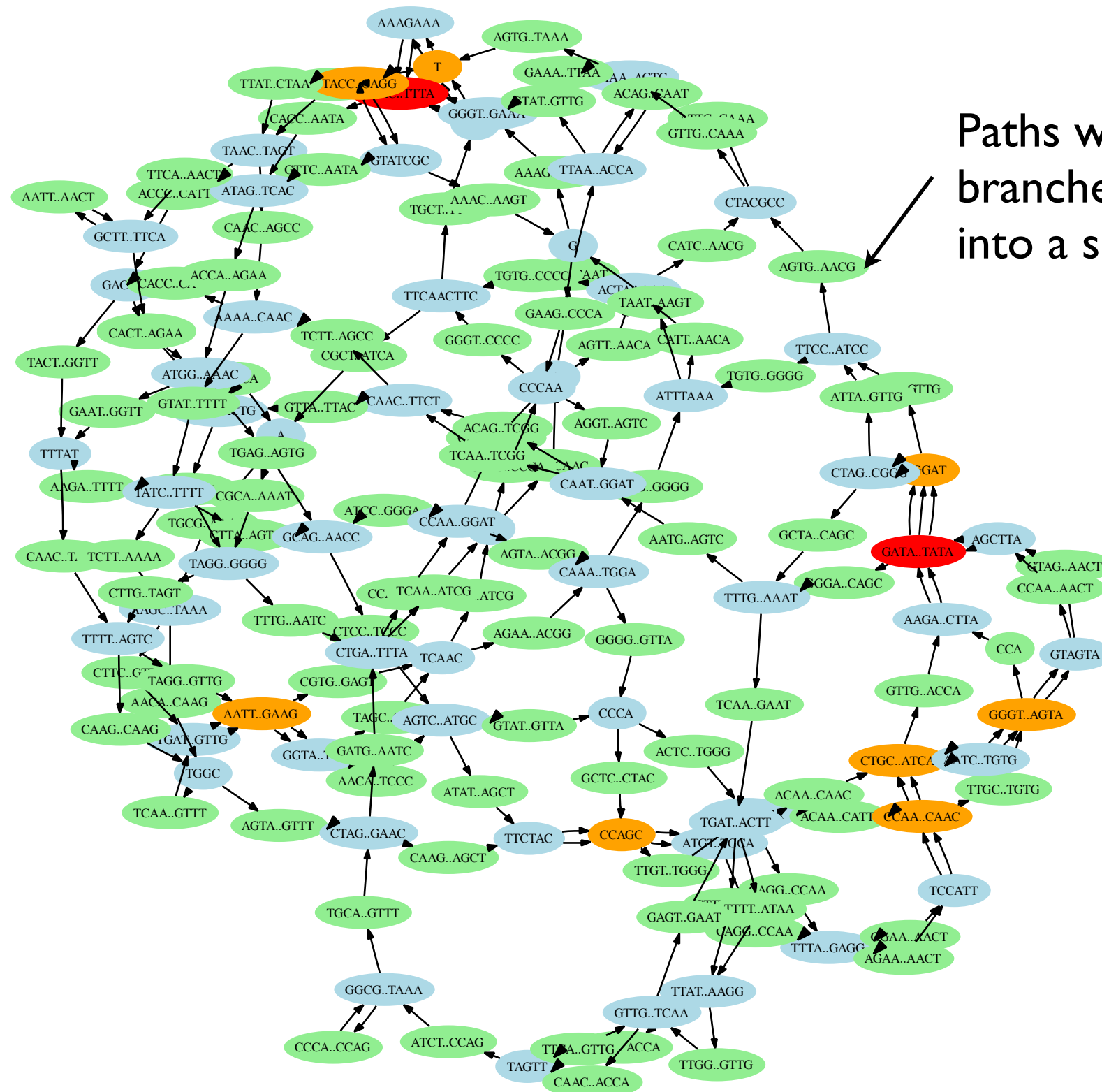


de Bruijn graph: nodes represent kmers, edges connect k-mers that are known to follow each other based on an observed read.

Can have > 1 edge between nodes.



Example bacterial de Bruijn graph

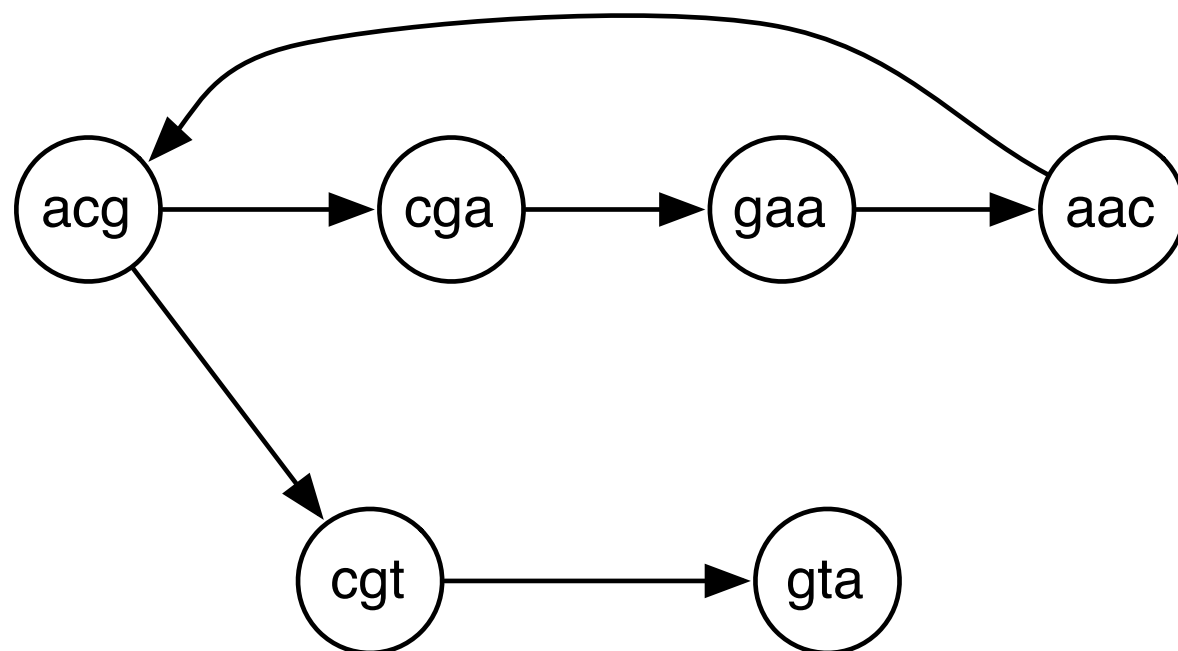


Paths with no branches compressed into a single node

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



acgaacgta

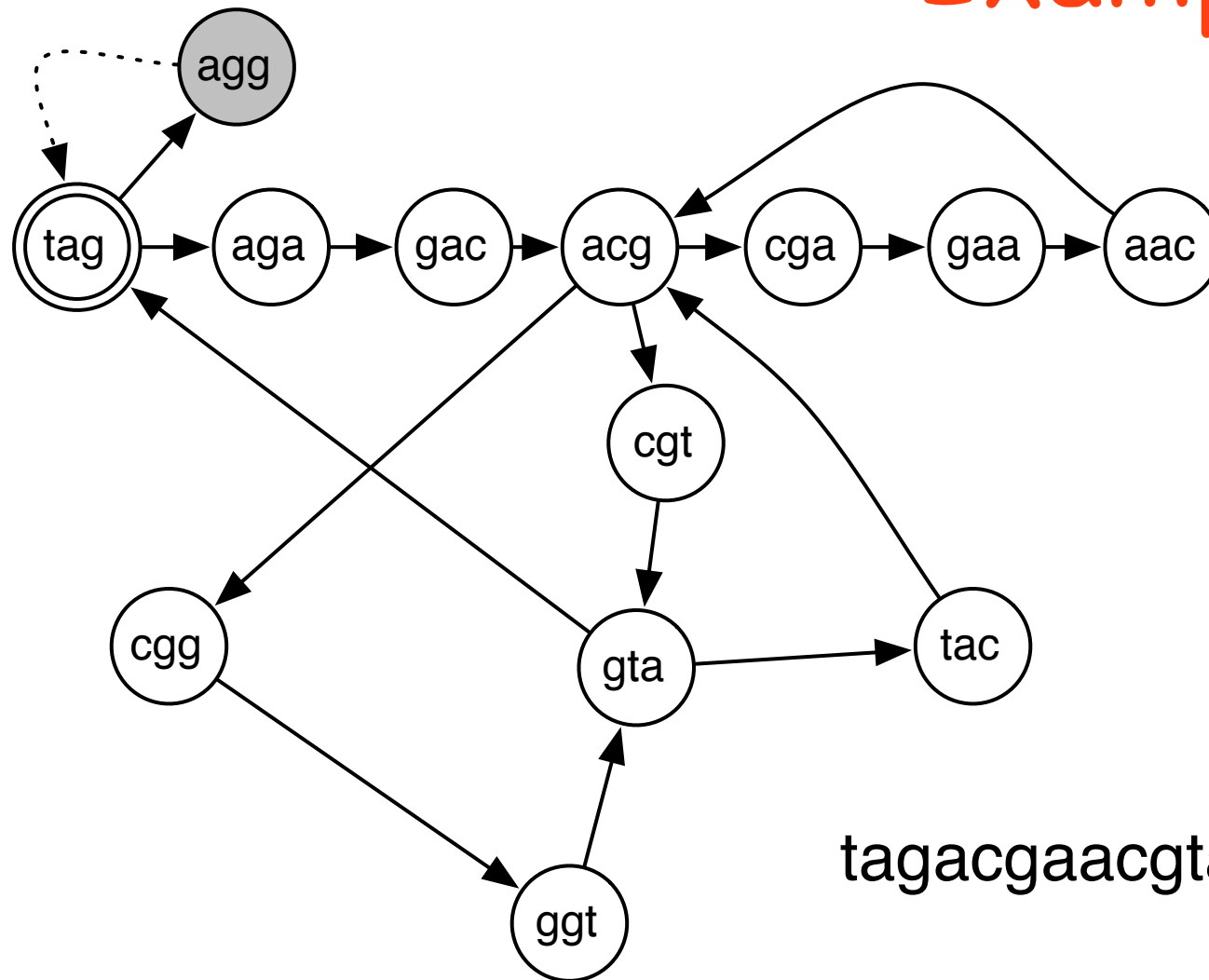
Let $dG(s)$ be the de Bruijn graph of string s . Then s corresponds to some Eulerian path in $dG(s)$.

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?

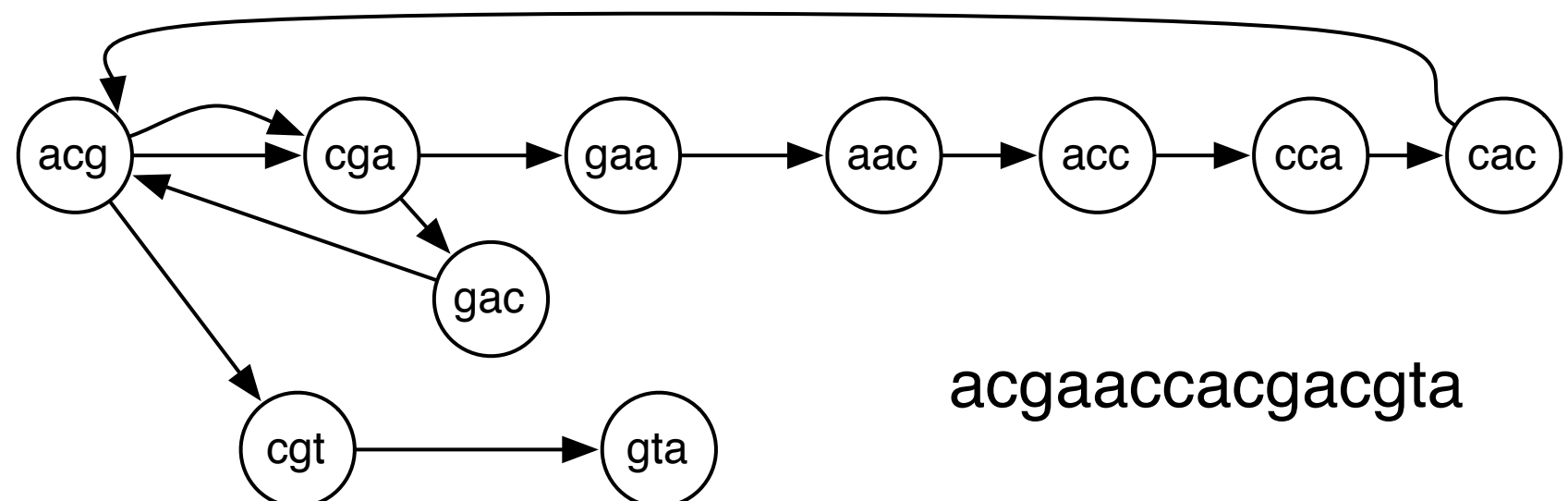
Examples



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

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

tagacgaacgtagcggtagg



acgaaccacgacgta

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.

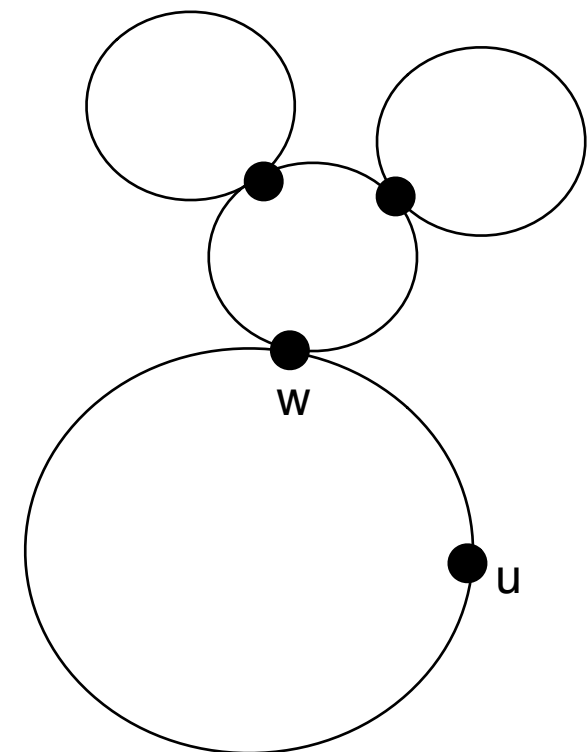
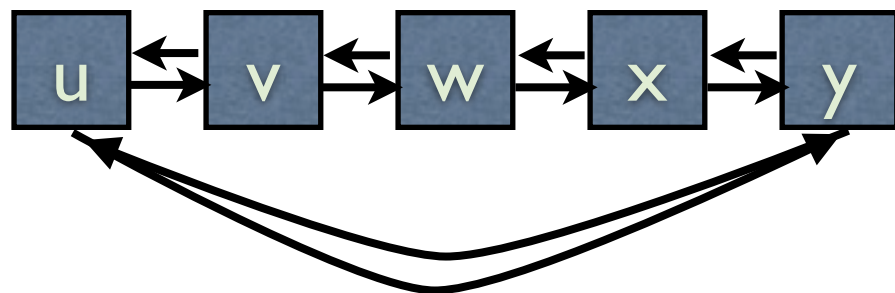
Why will you return to u ?

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

Repeat until all edges used:

- Start from some node w on the current tour with unused edges*.
- Walk along unused edges until you return to w , inserting the visited nodes after w into the current tour list.

*How can find such a node quickly?



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.

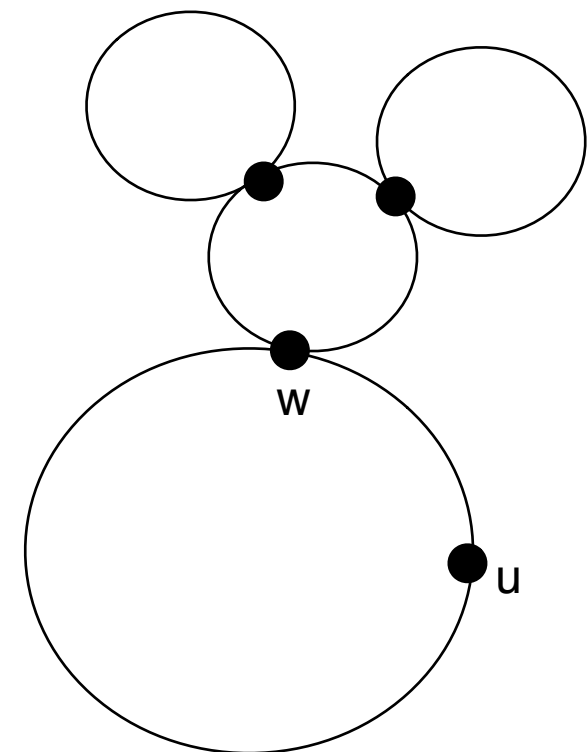
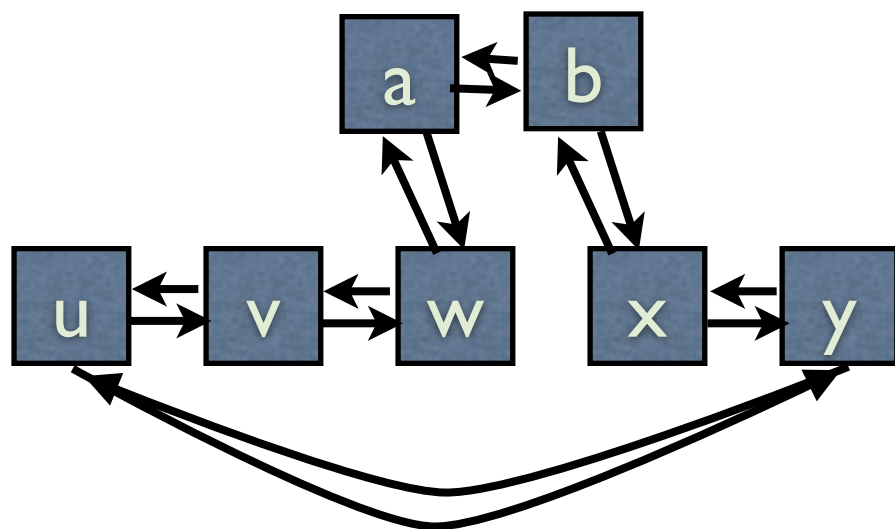
Why will you return to u ?

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

Repeat until all edges used:

- Start from some node w on the current tour with unused edges*.
- Walk along unused edges until you return to w , inserting the visited nodes after w into the current tour list.

*How can find such a node quickly?

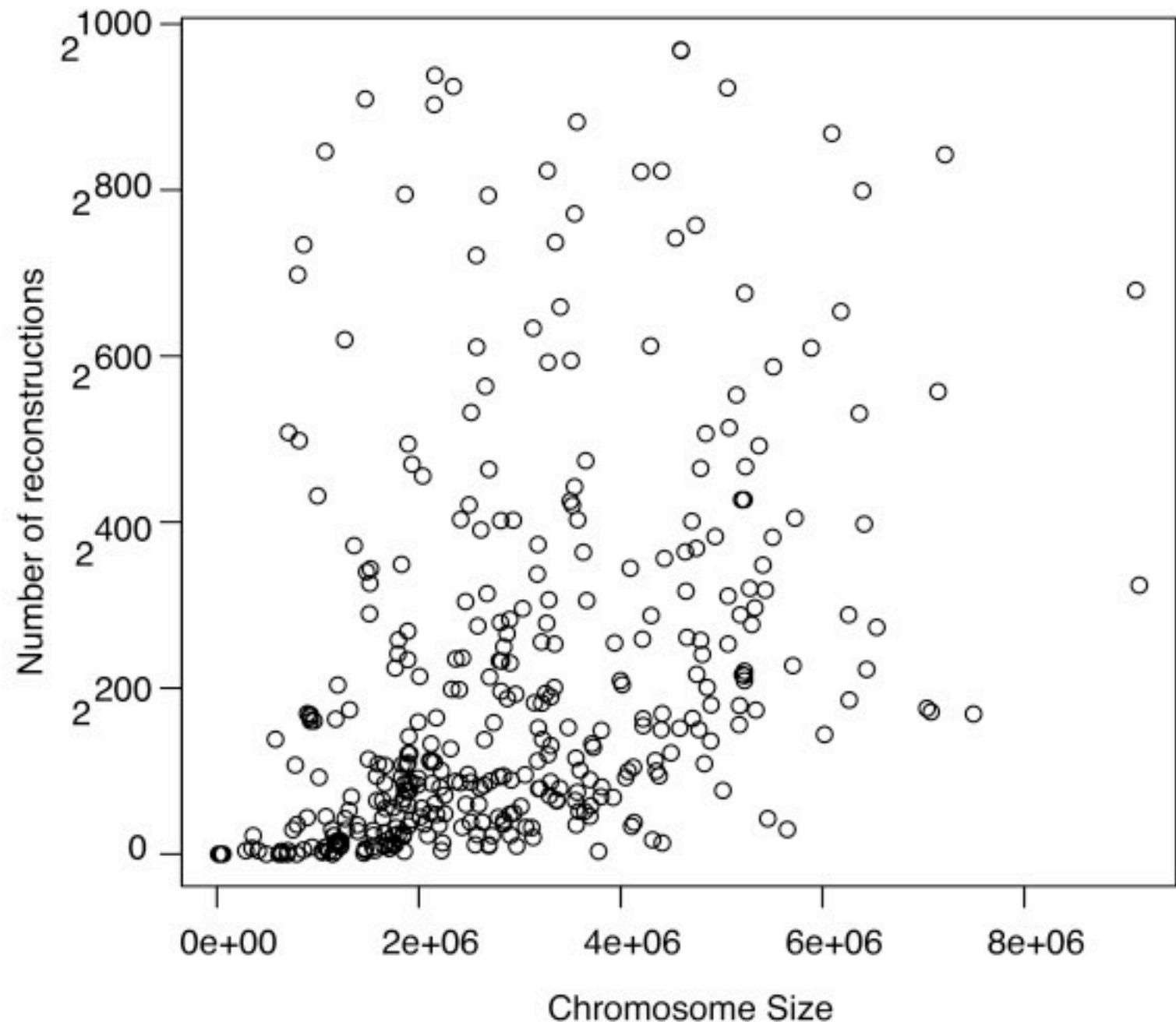


The Problem with Eulerian Paths

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.

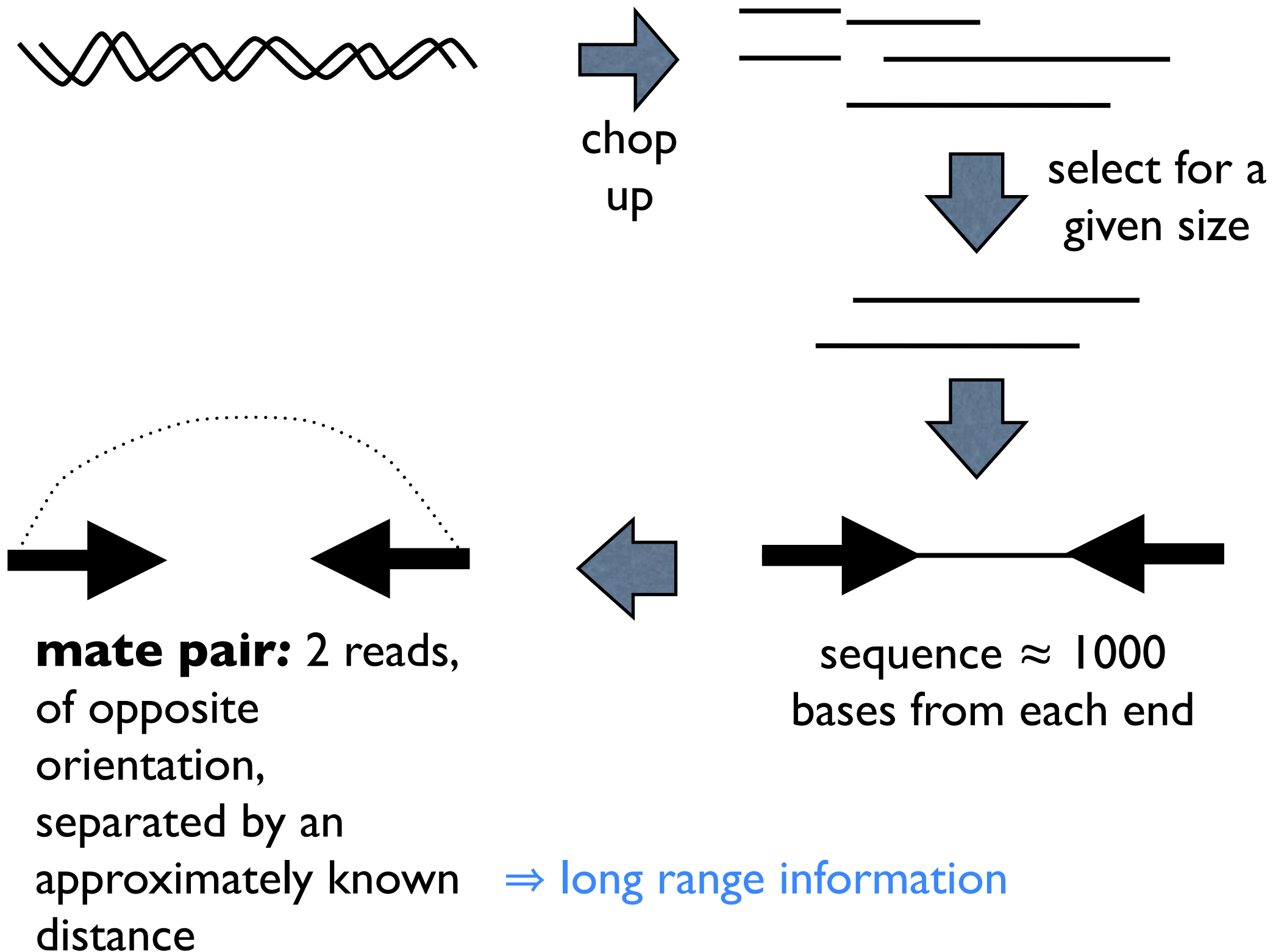
With imperfect data, there are usually NO Eulerian tours.



(Kingsford, Schatz, Pop, 2010)

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

Mate Pairs



References

- http://www.cbcb.umd.edu/research/assembly_primer
- <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2874646/>
- http://www.math.ucsd.edu/~gptesler/186/slides/shotgun_f13-handout.pdf
- <http://www.biomedcentral.com/1471-2105/11/21/abstract>