

Sequence Assembly

Fall 2016

BMI/CS 576

www.biostat.wisc.edu/bmi576/

Irene Ong

Irene.ong@wisc.edu

Outline

- What Is Genome Sequencing?
- Exploding Newspapers
- The String Reconstruction Problem
- **String Reconstruction as a Hamiltonian Path Problem**
- String Reconstruction as an Eulerian Path Problem
- Similar Problems with Different Fates
- De Bruijn Graphs
- Euler's Theorem
- Assembling Read-Pairs
- De Bruijn Graphs Face Harsh Realities of Assembly

Representing a Genome as a Path

Composition₃(TAATGCCATGGGATGTT) =

TAA AAT ATG TGC GCC CCA CAT ATG TGG GGG GGA GAT ATG TGT GTT

Representing a Genome as a Path

Composition₃(TAATGCCATGGGATGTT) =



Representing a Genome as a Path

Composition₃(TAATGCCATGGGATGTT) =



Can we construct this **genome path** without knowing the genome
TAATGCCATGGGATGTT, only from its composition?

Representing a Genome as a Path

Composition₃(TAATGCCATGGGATGTT) =

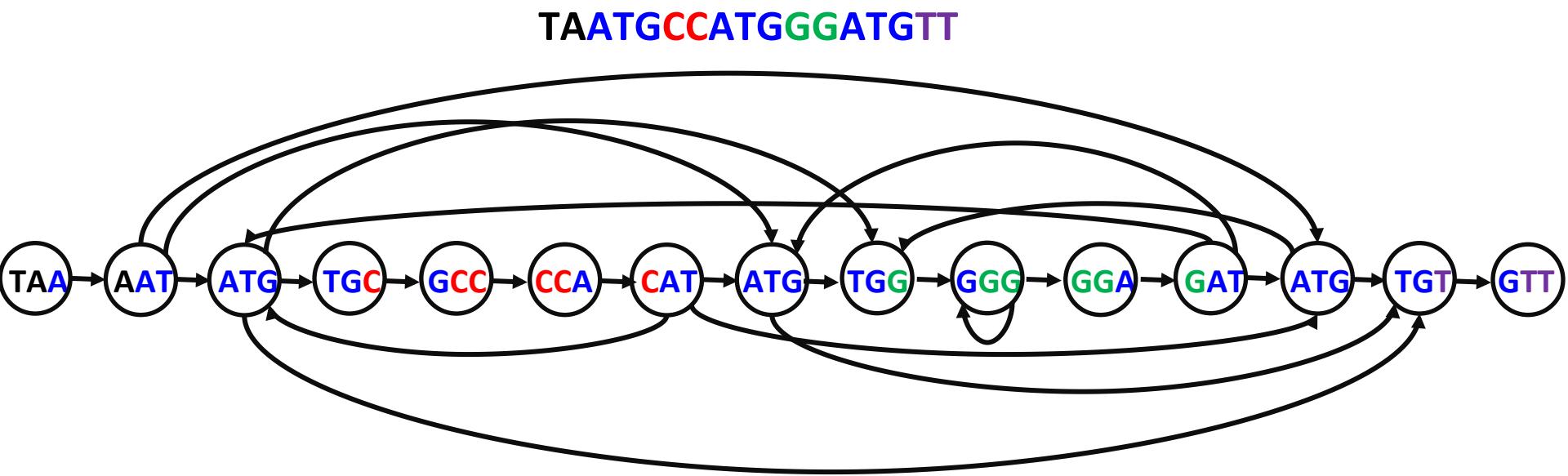


Can we construct this **genome path** without knowing the genome
TAATGCCATGGGATGTT, only from its composition?

Yes. We simply need to connect k-mer₁ with k-mer₂ if suffix(k-mer₁)=prefix(k-mer₂).

E.g. TAA → AAT

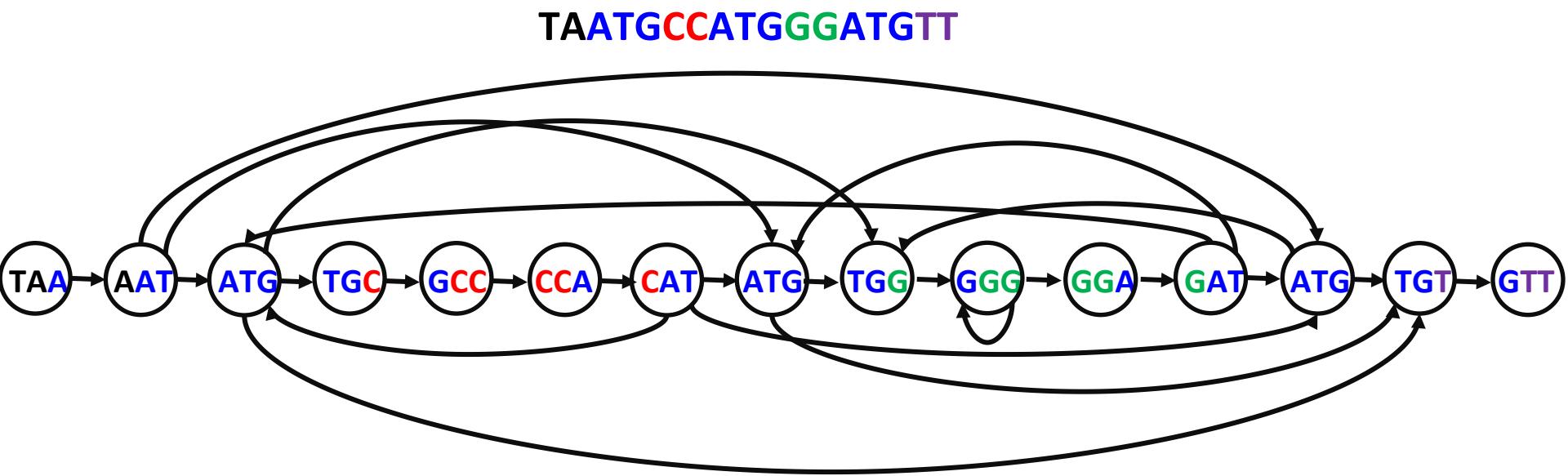
A Path Turns into a Graph



Yes. We simply need to connect $k\text{-mer}_1$ with $k\text{-mer}_2$ if $\text{suffix}(k\text{-mer}_1) = \text{prefix}(k\text{-mer}_2)$.

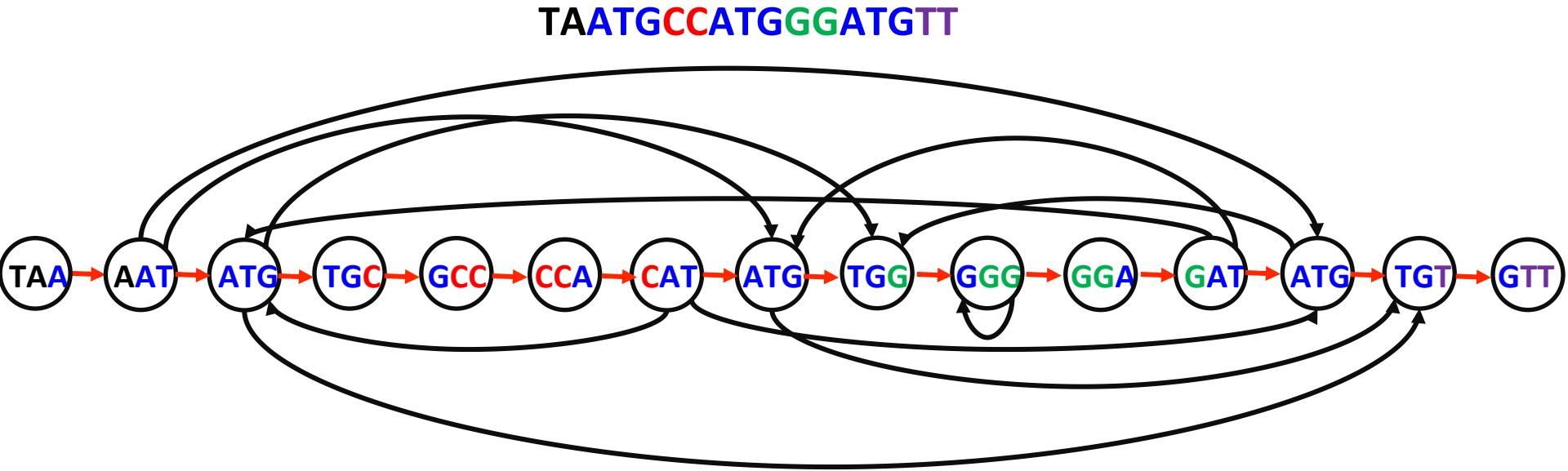
E.g. TAA \rightarrow AAT

A Path Turns into a Graph



Can we still find the **genome path** in this graph?

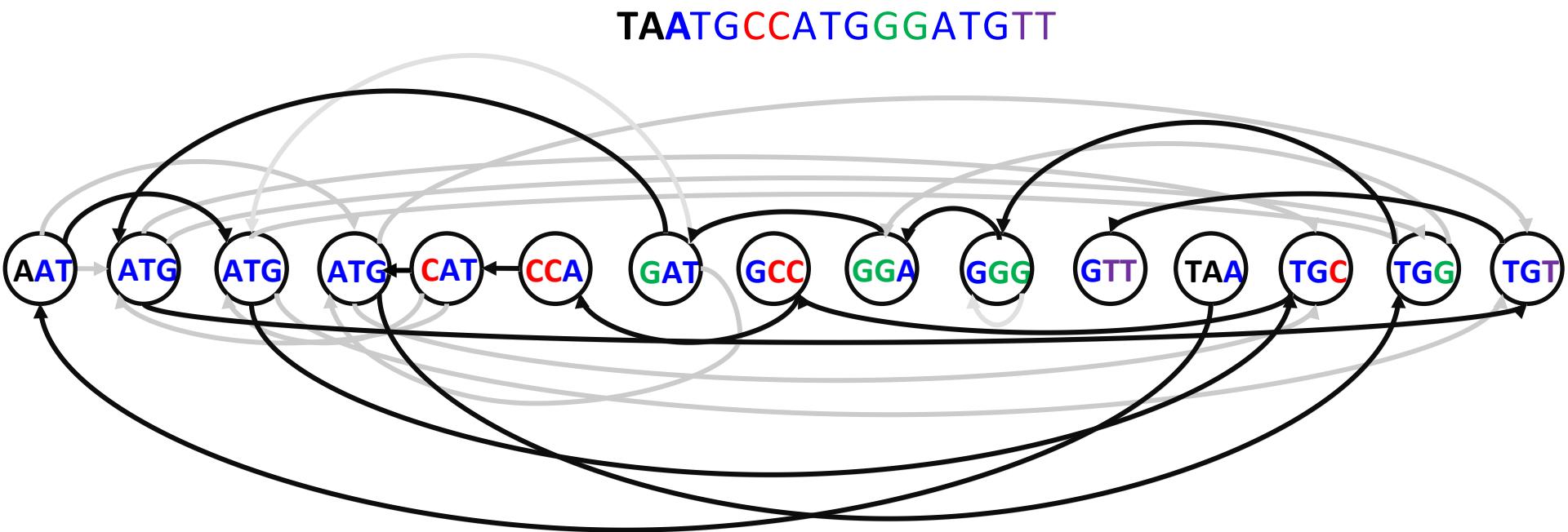
A Path Turns into a Graph



Can we still find the genome path in this graph?

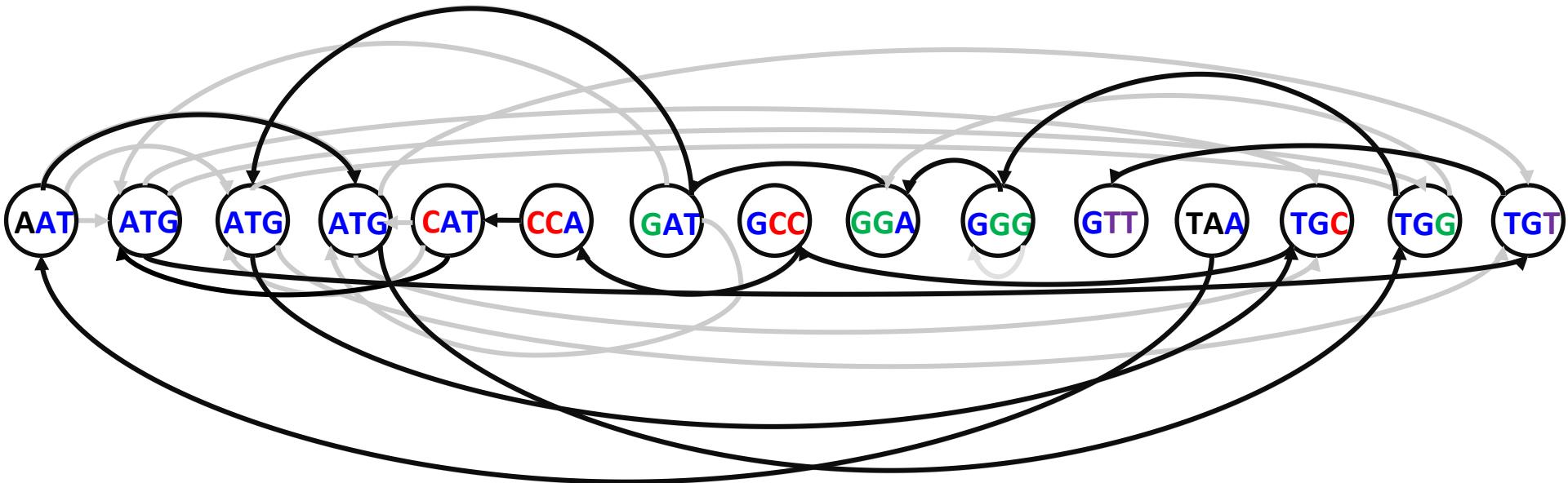
Where Is the Genomic Path?

A Hamiltonian path: a path that visits each node in a graph exactly once.

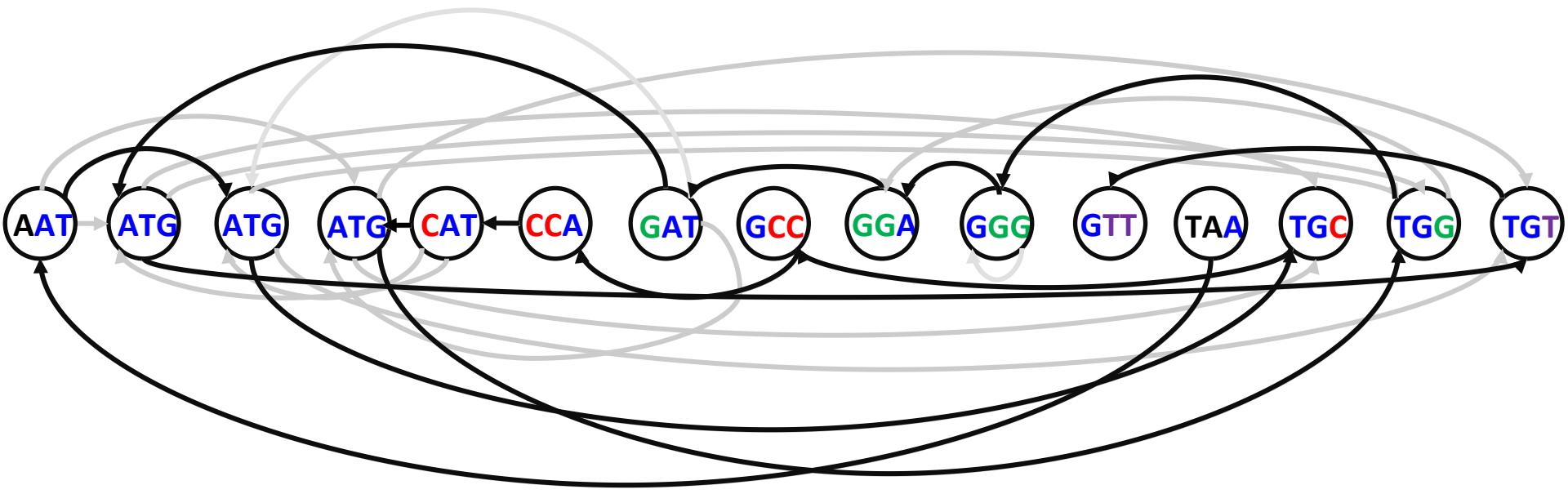


What are we trying to find in this graph?

TAATGGGATG CCATGTT



TAATGCCATGGGATGTT



Outline

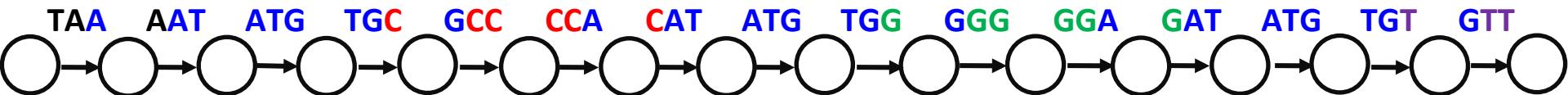
- What Is Genome Sequencing?
- Exploding Newspapers
- The String Reconstruction Problem
- String Reconstruction as a Hamiltonian Path Problem
- **String Reconstruction as an Eulerian Path Problem**
- Similar Problems with Different Fates
- De Bruijn Graphs
- Euler's Theorem
- Assembling Read-Pairs
- De Bruijn Graphs Face Harsh Realities of Assembly

A Slightly Different Path

TAATGCCATGGGATGTT

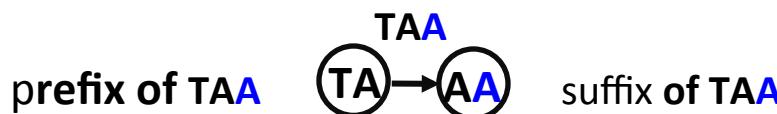


Instead of representing 3-mers as nodes



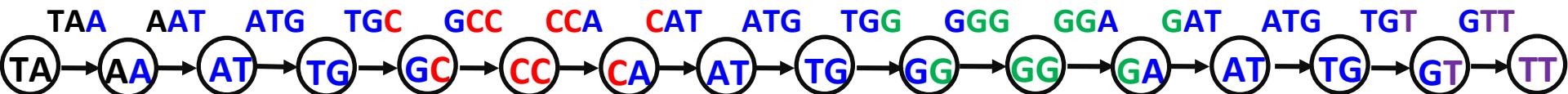
Let's represent 3-mers as edges

How do we label the starting and ending nodes of an edge?

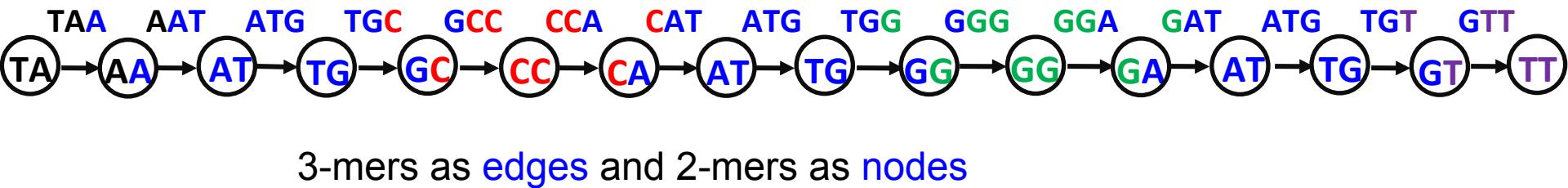


Labeling Nodes in the New Path

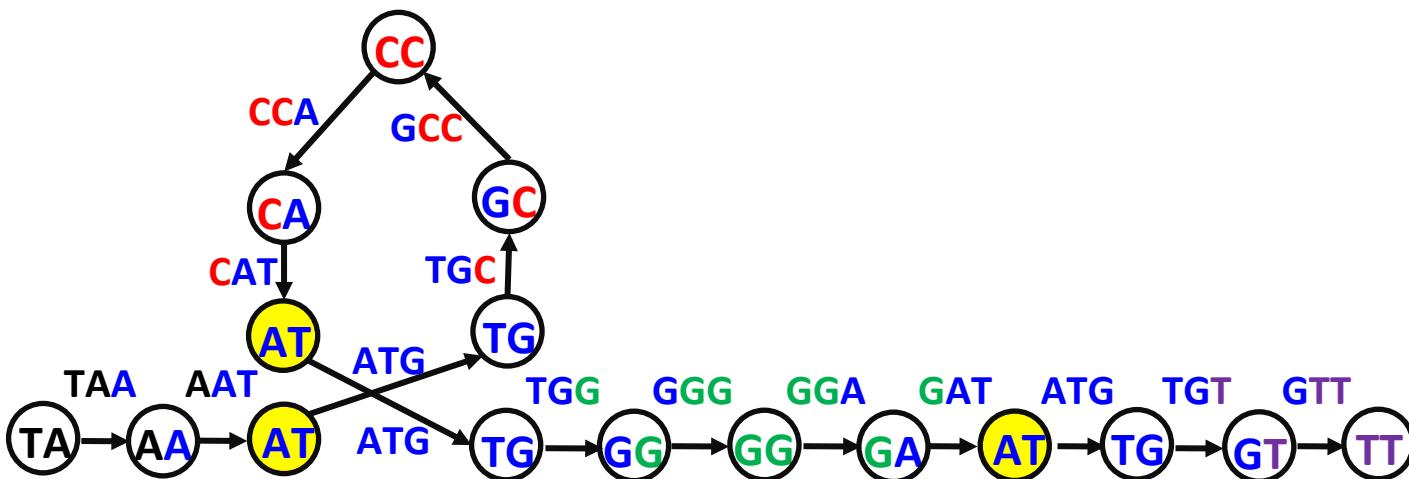
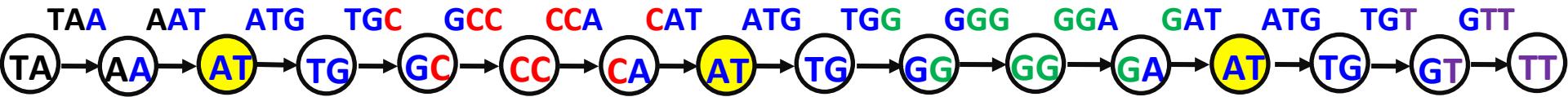
TAATGCCATGGGATGTT



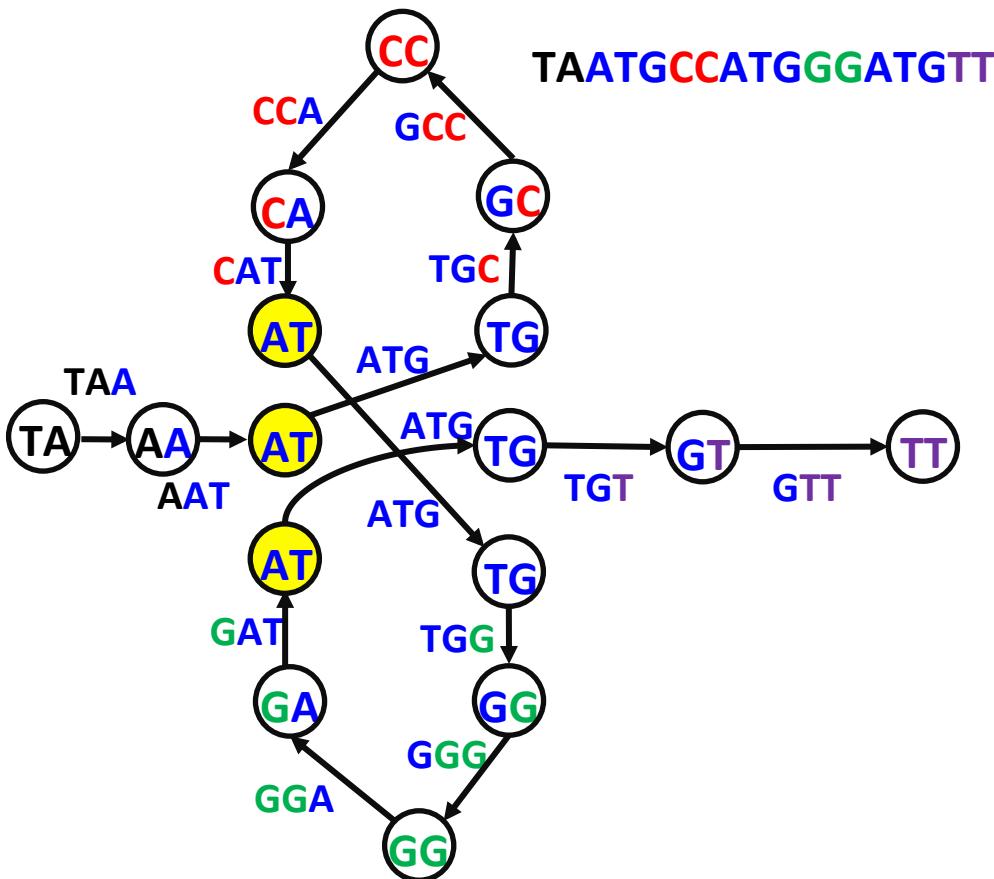
Labeling Nodes in the New Path



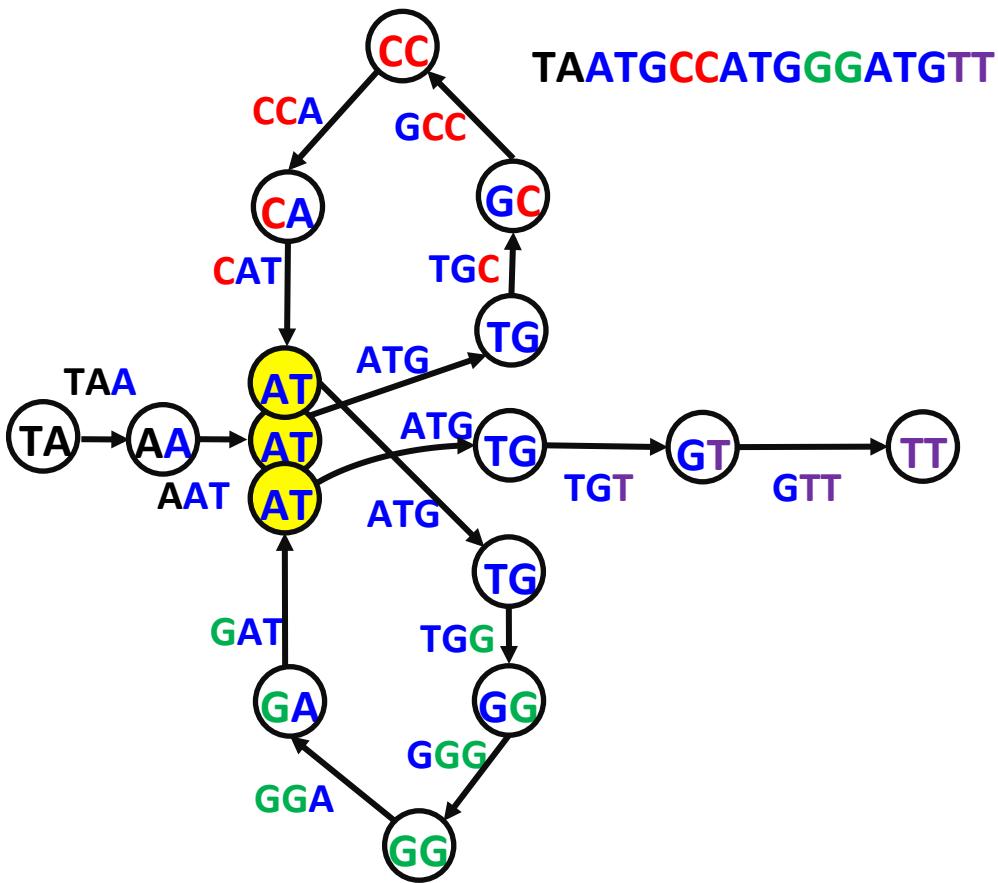
Gluing Identically Labeled Nodes



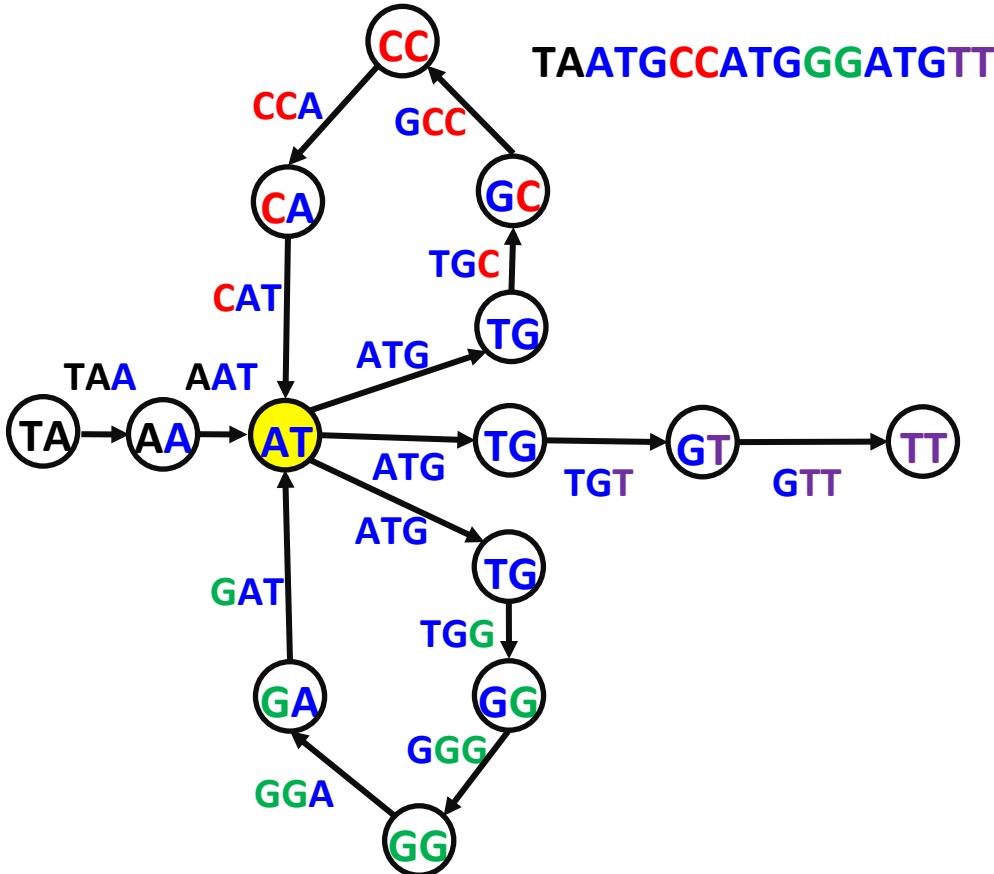
Gluing Identically Labeled Nodes



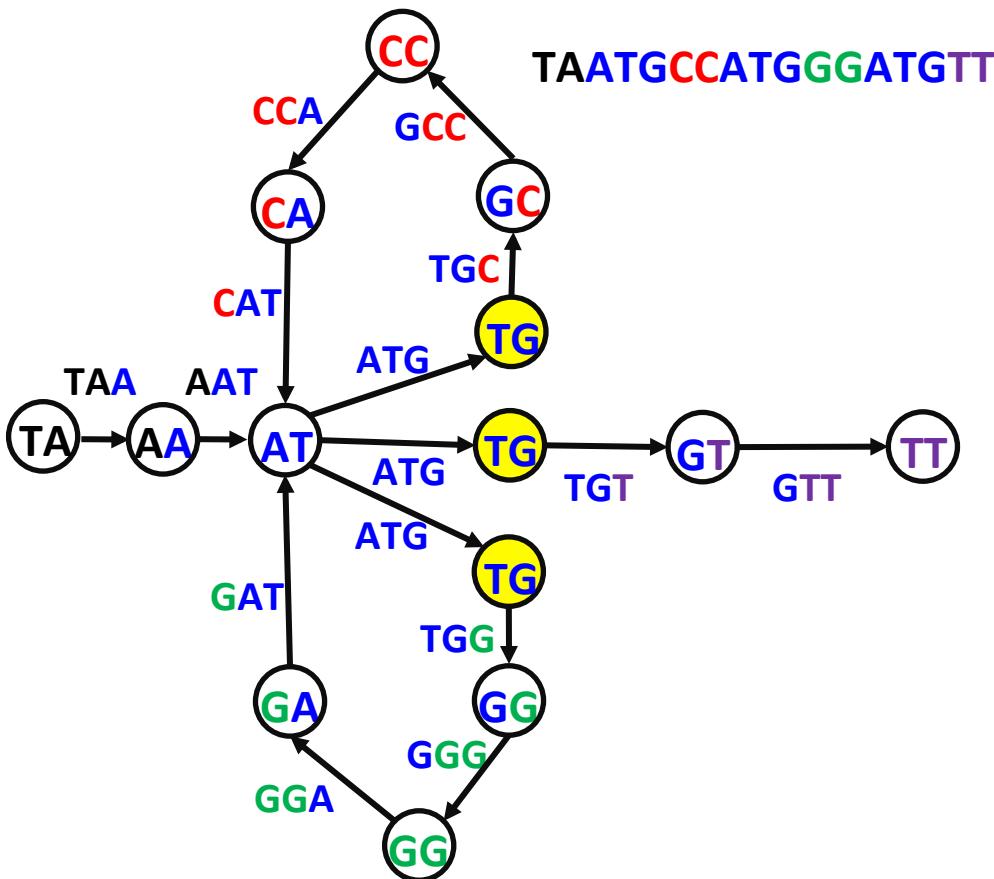
Gluing Identically Labeled Nodes



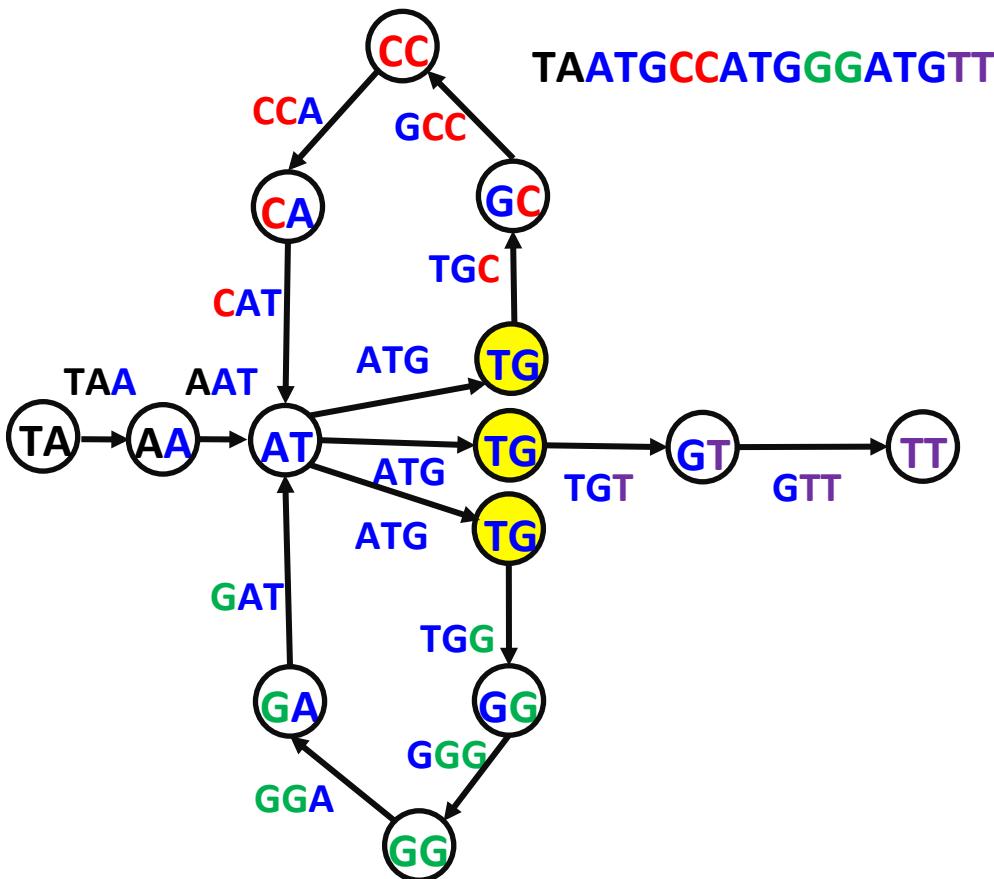
Gluing Identically Labeled Nodes



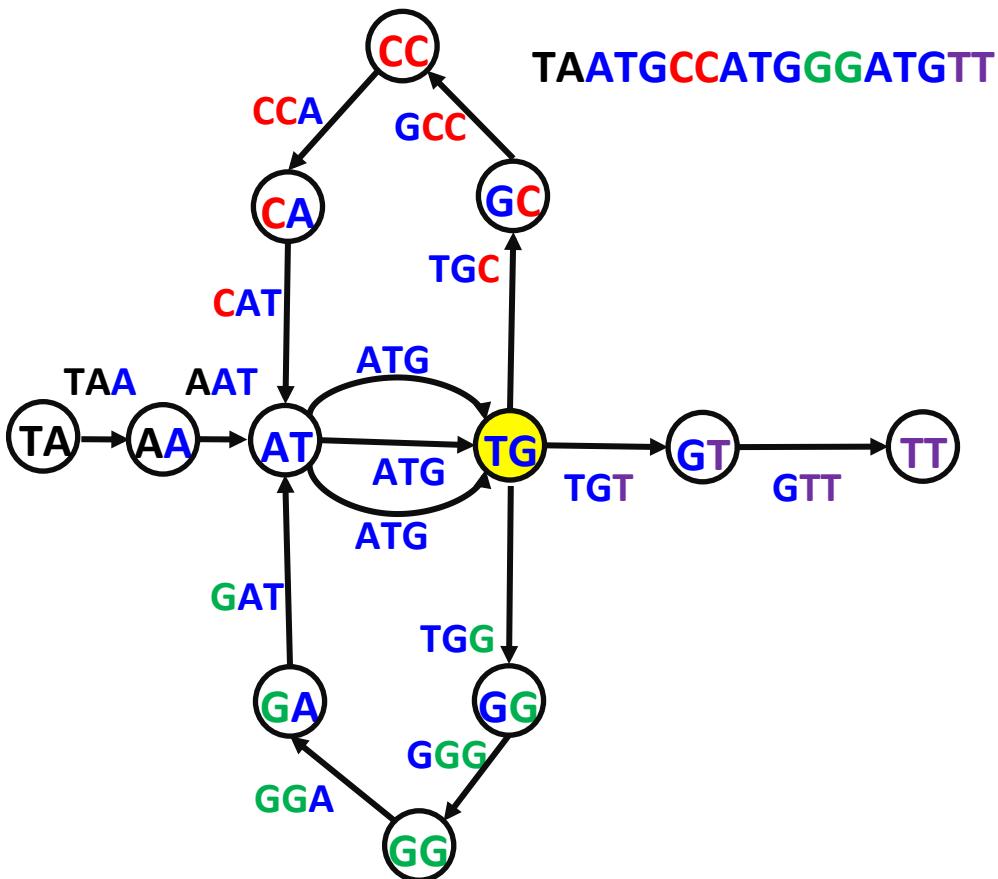
Gluing Identically Labeled Nodes



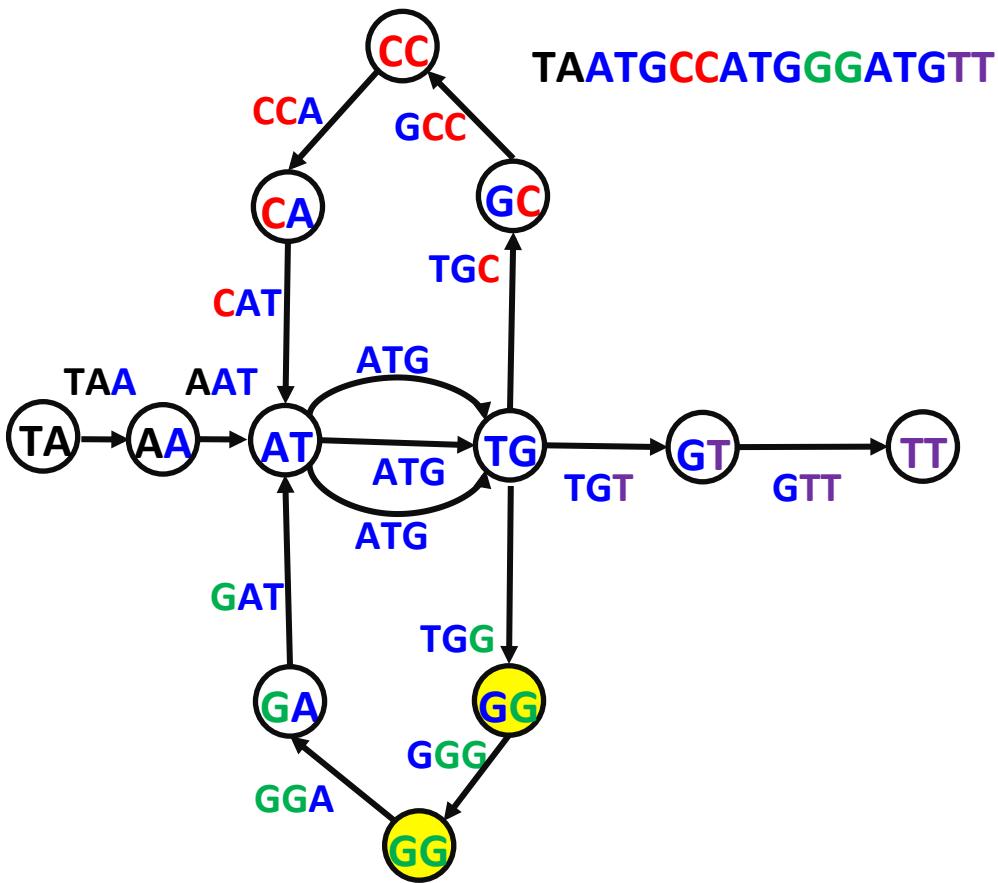
Gluing Identically Labeled Nodes



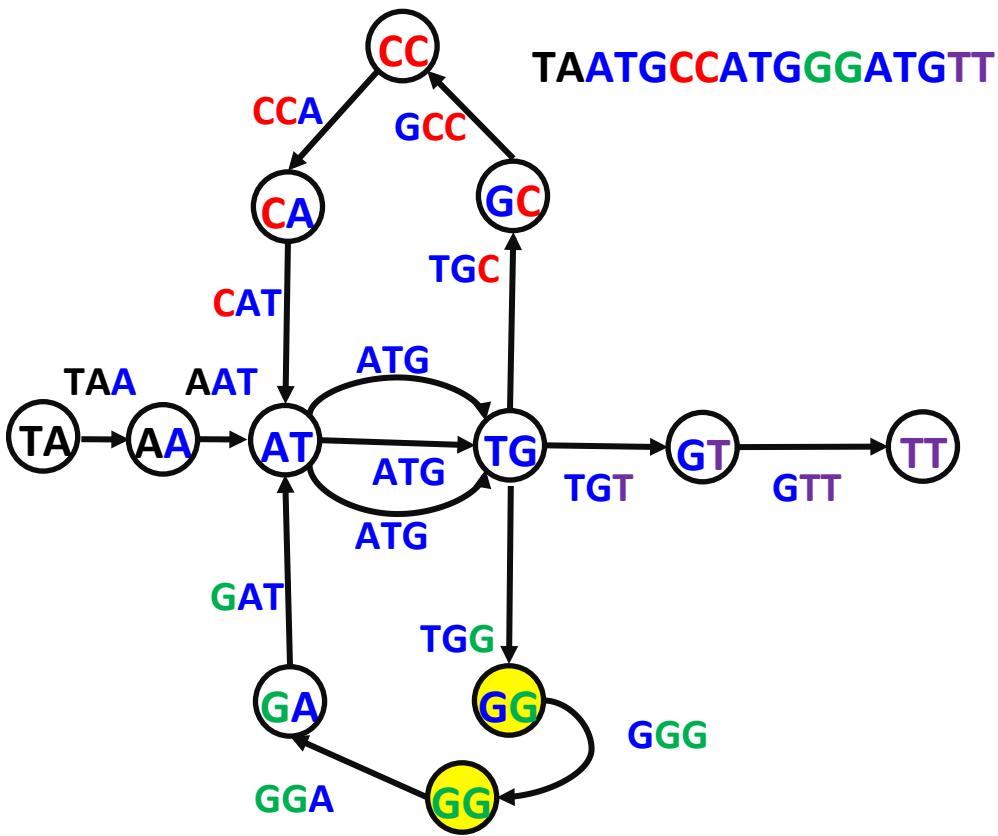
Gluing Identically Labeled Nodes



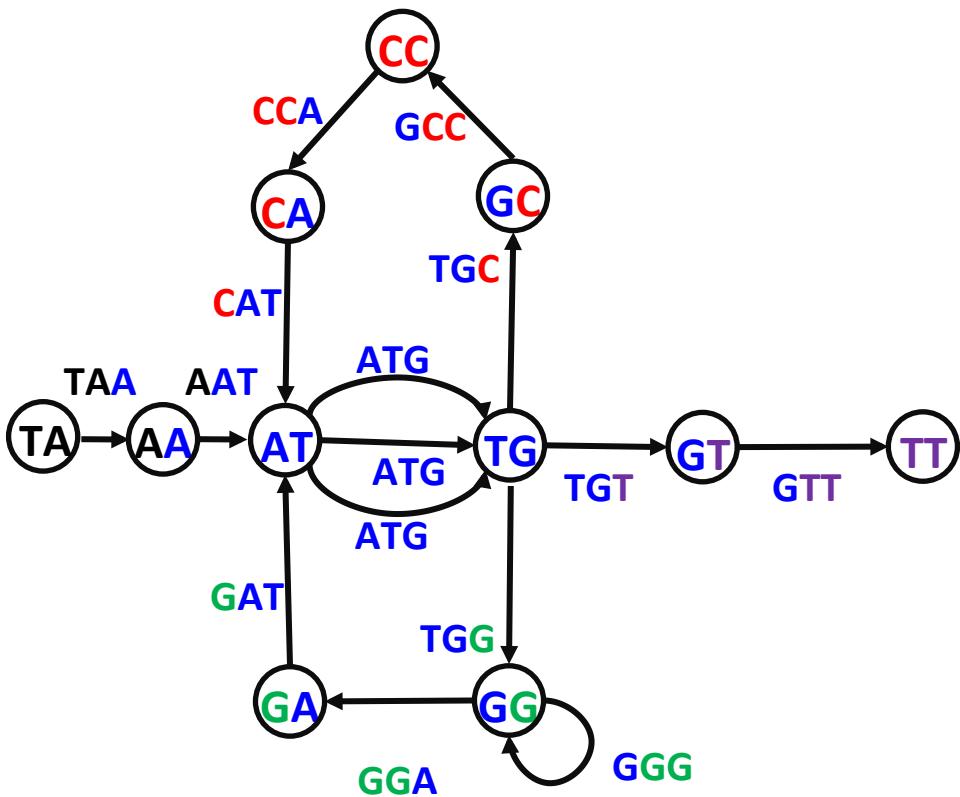
Gluing Identically Labeled Nodes



Gluing Identically Labeled Nodes



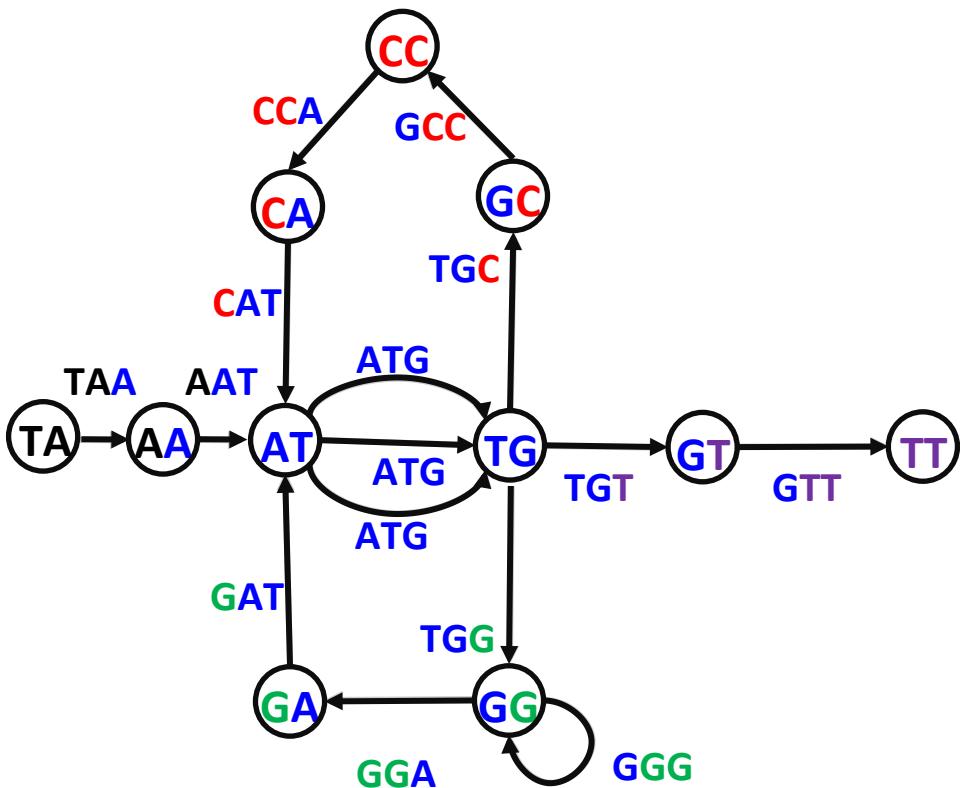
De Bruijn Graph of TAATGCCATGGGATGTT



Where is the Genome
hiding in this graph?

It Was Always There!

TAATGCCATGGGATGTT



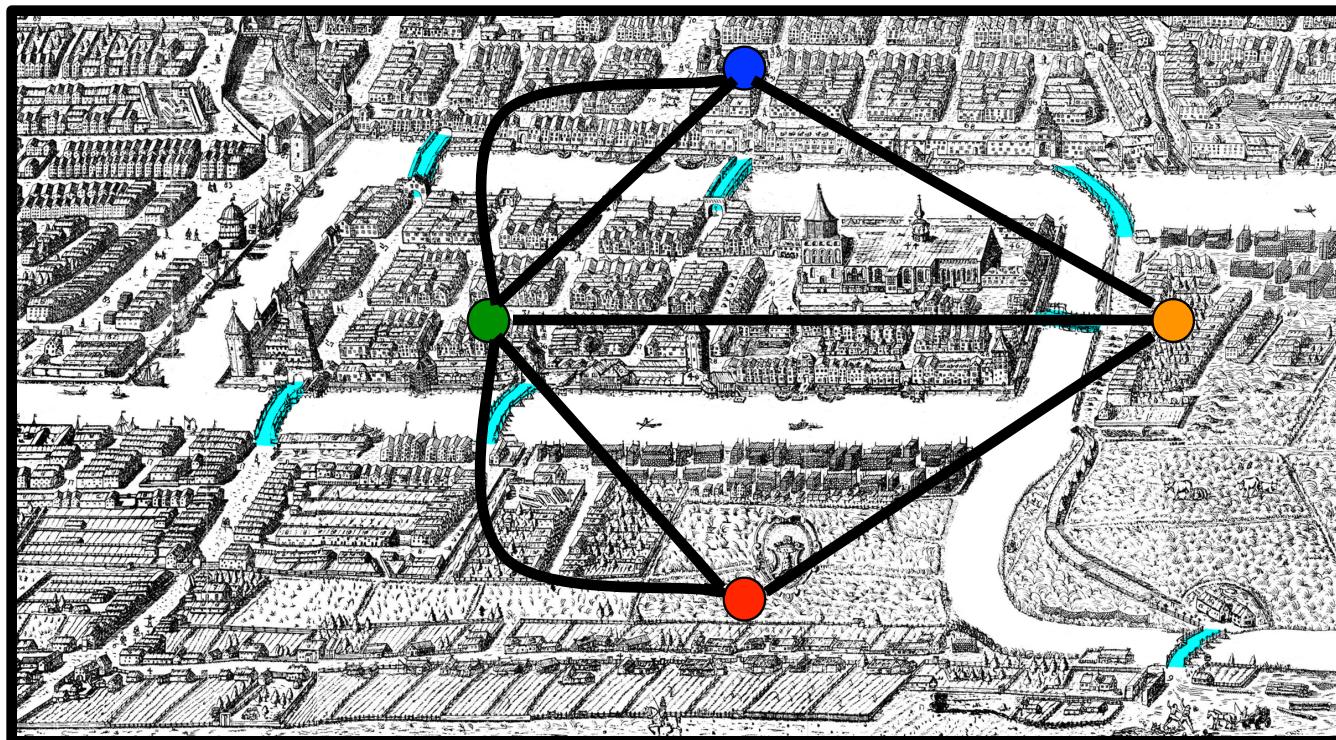
An Eulerian **path** in a graph is a path that visits each **edge** exactly once.

Eulerian Path Problem

Eulerian Path Problem. Find an Eulerian path in a graph.



- Input. A graph.
- Output. A path visiting every edge in the graph exactly once.



Outline

- What Is Genome Sequencing?
- Exploding Newspapers
- The String Reconstruction Problem
- String Reconstruction as a Hamiltonian Path Problem
- String Reconstruction as an Eulerian Path Problem
- **Similar Problems with Different Fates**
- De Bruijn Graphs
- Euler's Theorem
- Assembling Read-Pairs
- De Bruijn Graphs Face Harsh Realities of Assembly

Eulerian Versus Hamiltonian Paths

Eulerian Path Problem. Find an **Eulerian** path in a graph.



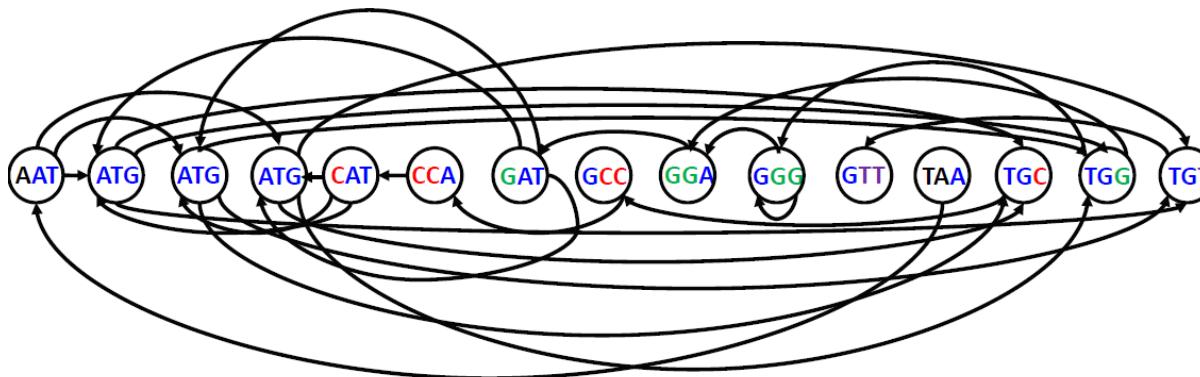
- Input. A graph.
- Output. A path visiting every **edge** in the graph exactly once.

Hamiltonian Path Problem. Find a **Hamiltonian** path in a graph.

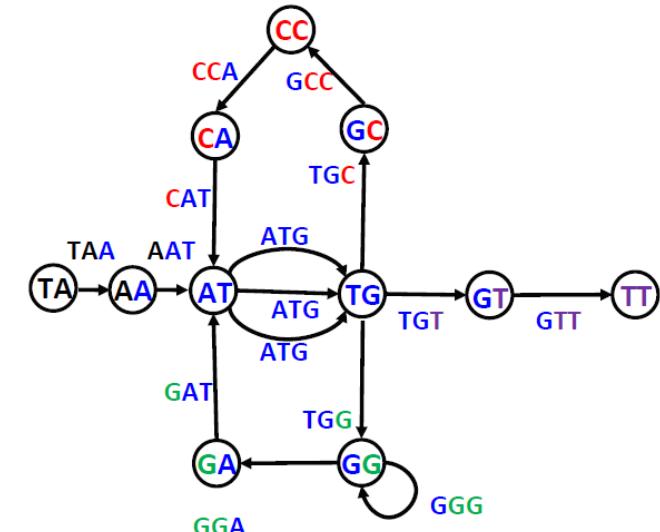


- Input. A graph.
- Output. A path visiting every **node** in the graph exactly once.

What Problem Would You Prefer to Solve?

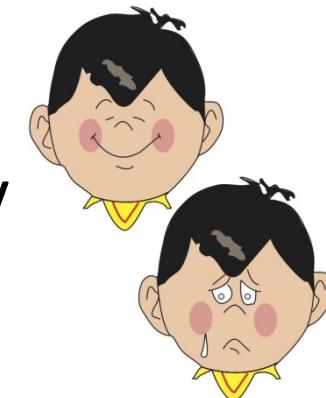


Hamiltonian Path Problem



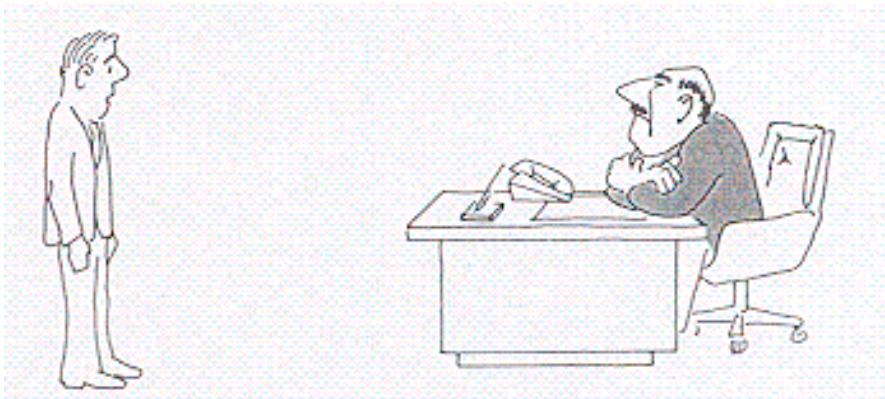
Eulerian Path Problem

While Euler solved the Eulerian Path Problem (even for a city with a million bridges), nobody has developed a fast algorithm for the Hamiltonian Path Problem yet.



NP-Complete Problems

- The Hamiltonian Path Problem belongs to a collection containing thousands of computational problems for which no fast algorithms are known.



“I can't find an efficient algorithm, I must not be working hard enough.”

Change of Attitude

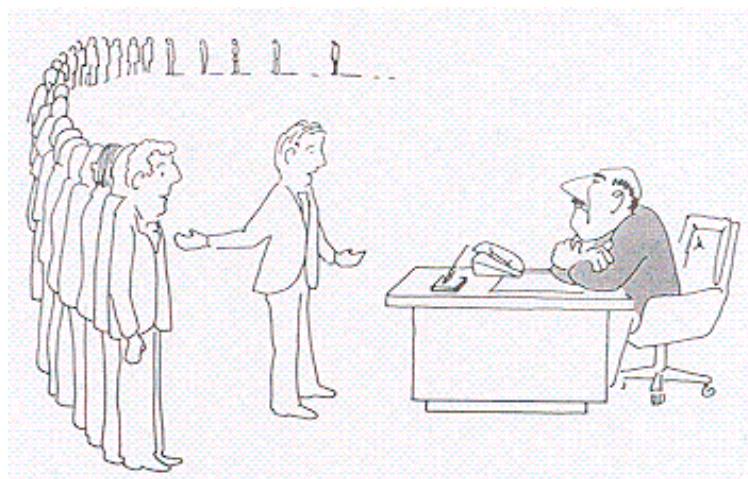
That would be an excellent argument, but the question of whether or not NP-Complete problems can be solved efficiently is one of seven **Millennium Problems** in mathematics.



“I can't find an efficient algorithm, because no such algorithm is possible.”

The Modern State of Affairs

NP-Complete problems are all equivalent: find an efficient solution to one, and you have an efficient solution to them all.



“I can't find an efficient algorithm, but neither can all these famous people.”

Outline

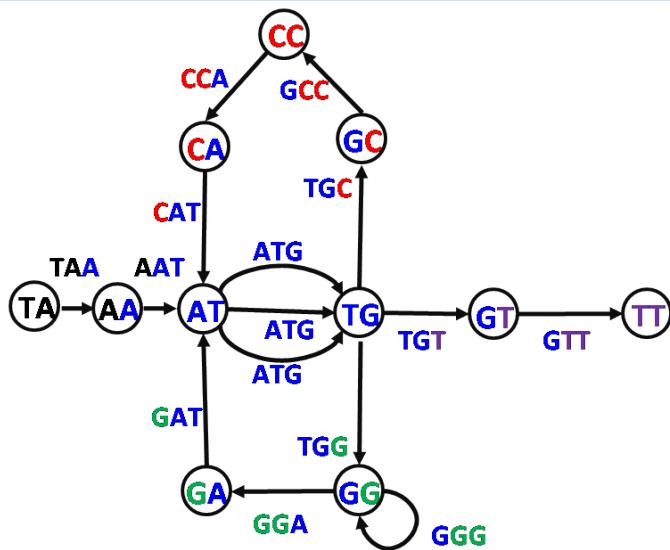
- What Is Genome Sequencing?
- Exploding Newspapers
- The String Reconstruction Problem
- String Reconstruction as a Hamiltonian Path Problem
- String Reconstruction as an Eulerian Path Problem
- Similar Problems with Different Fates
- **De Bruijn Graphs**
- Euler's Theorem
- Assembling Read-Pairs
- De Bruijn Graphs Face Harsh Realities of Assembly

Eulerian Path Problem

Eulerian Path Problem. Find an **Eulerian** path in a graph.

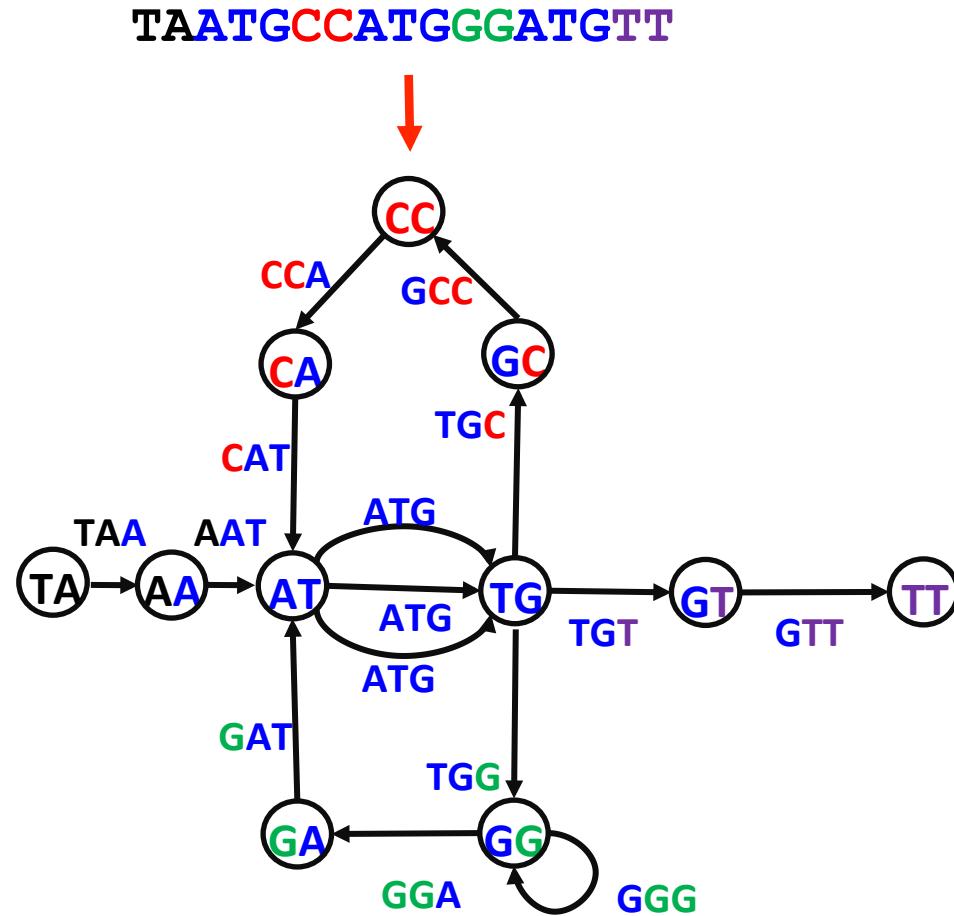


- Input. A graph.
- Output. A path visiting every **edge** in the graph exactly once.



We constructed the de Bruijn graph from Genome, but in reality, Genome is unknown!

What We Have Done: From Genome to de Bruijn Graph



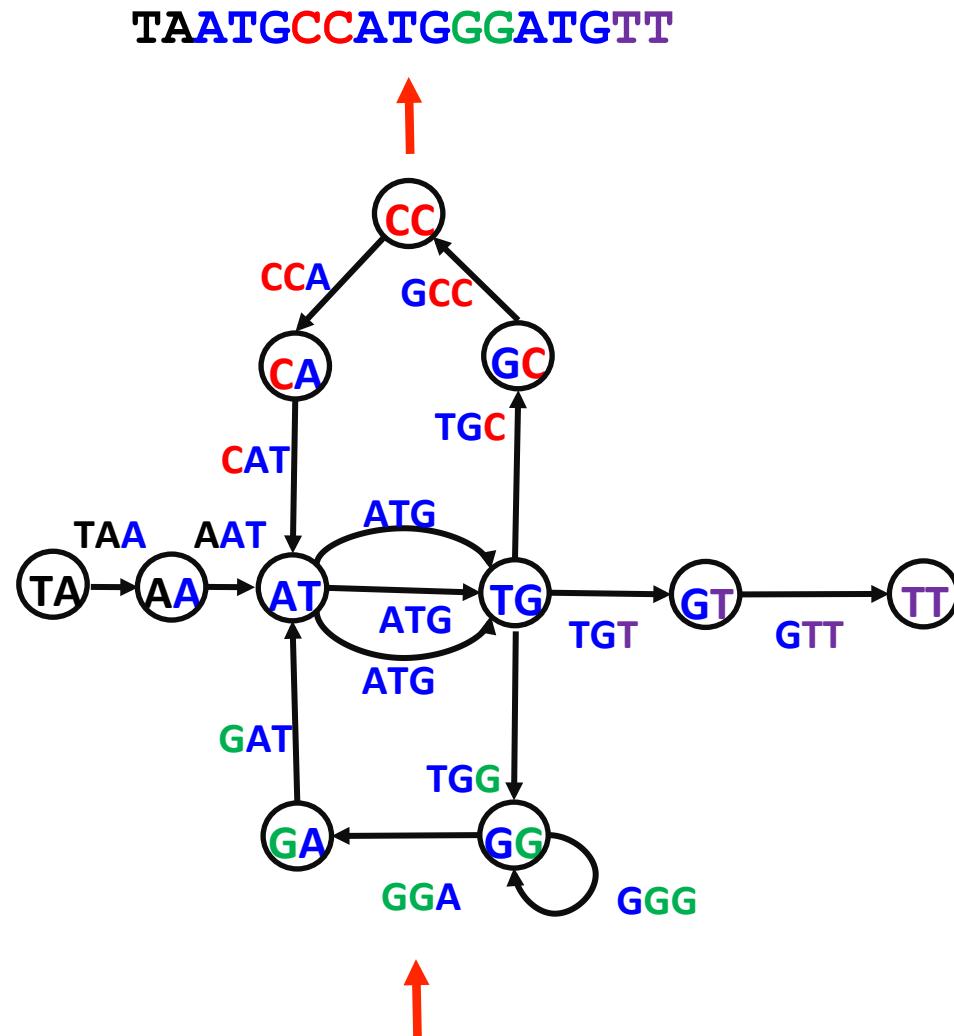
What We Want: From Reads (k-mers) to Genome

TAATGCCATGGGATGTT



AAT ATG ATG ATG CAT CCA GAT GCC GGA GGG GTT TAA TGC TGG TGT

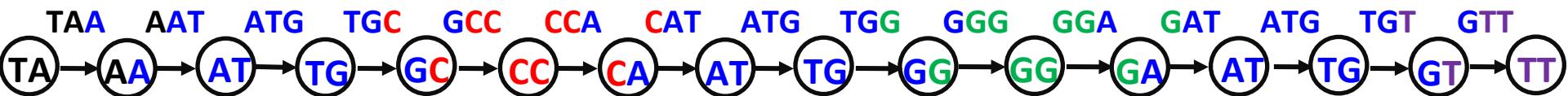
What We will Show: From Reads to de Bruijn Graph to Genome



AAT ATG ATG ATG CAT CCA GAT GCC GGA GGG GTT TAA TGC TGG TGT

Constructing de Bruijn Graph when Genome Is Known

TAATGCCATGGGATGTT



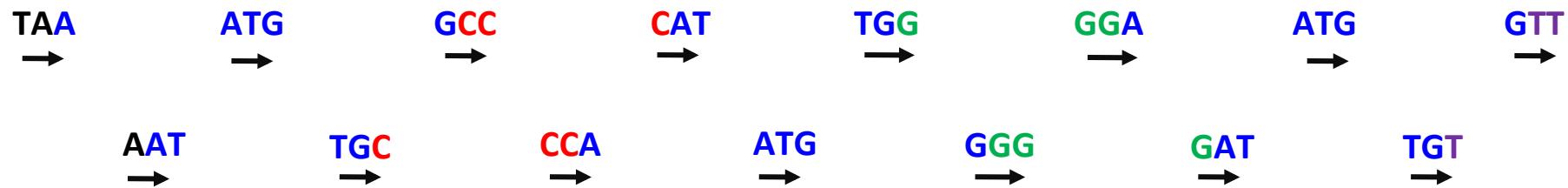
Constructing de Bruijn when Genome Is Unknown

TAA ATG GCC CAT TGG GGA ATG GTT

AAT TGC CCA ATG GGG GAT TGT

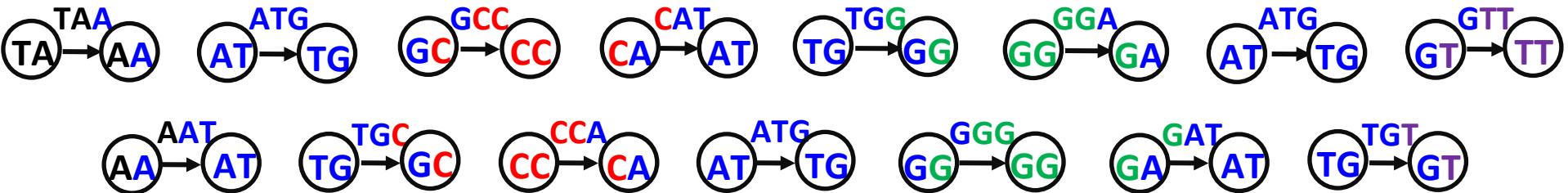
Composition₃(TAATGCCATGGGATGTT)

Representing Composition as a Graph Consisting of Isolated Edges



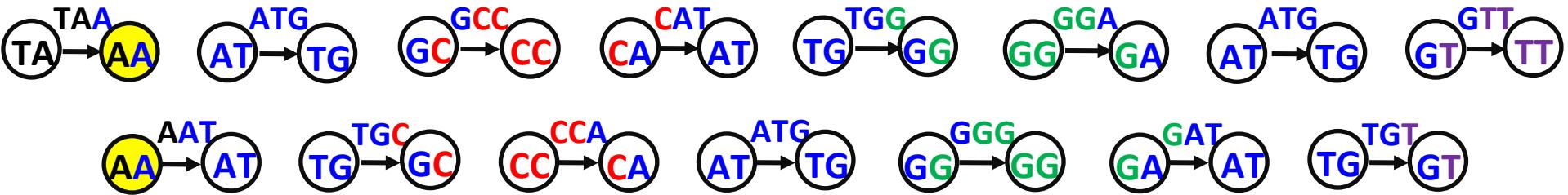
Composition₃(TAATGCCATGGGATGTT)

Constructing de Bruijn Graph from k-mer Composition

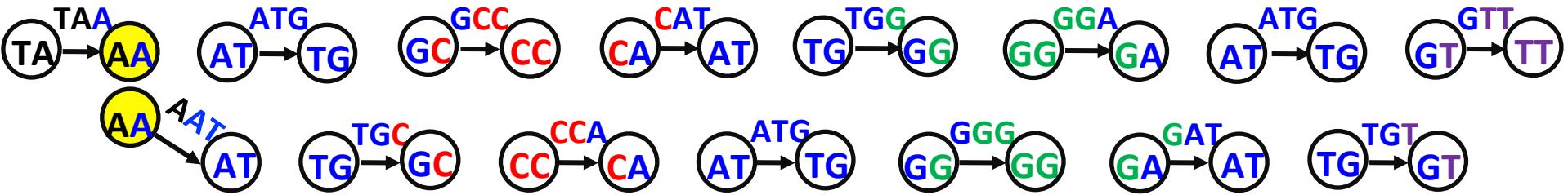


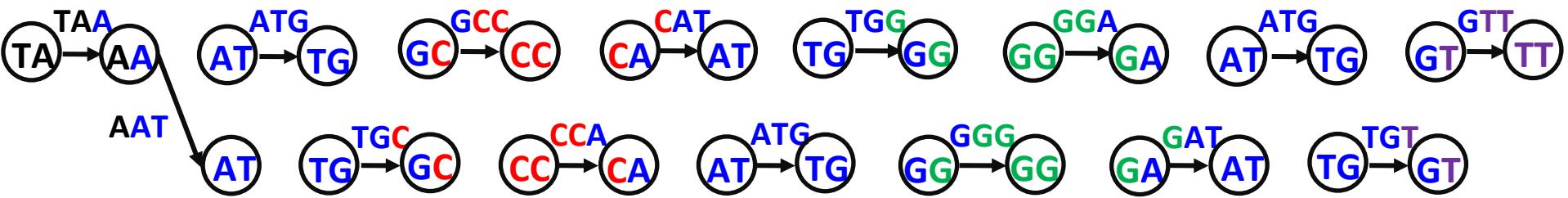
Composition₃(TAATGCCATGGGATGTT)

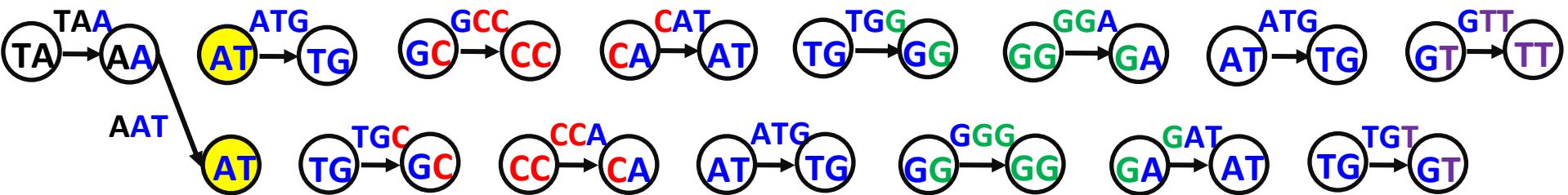
Gluing Identically Labeled Nodes

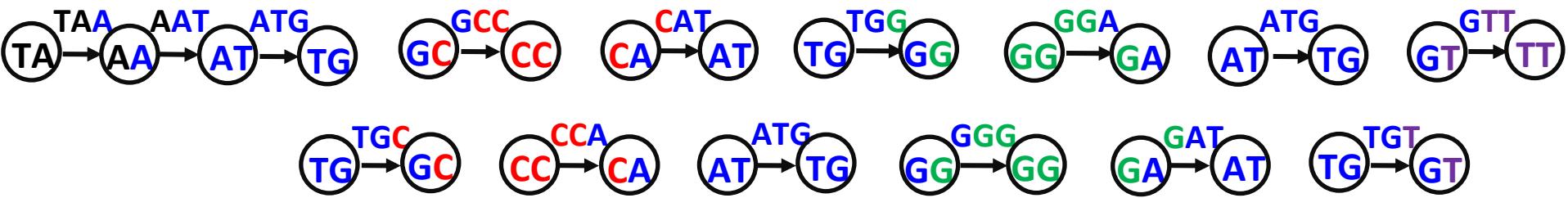


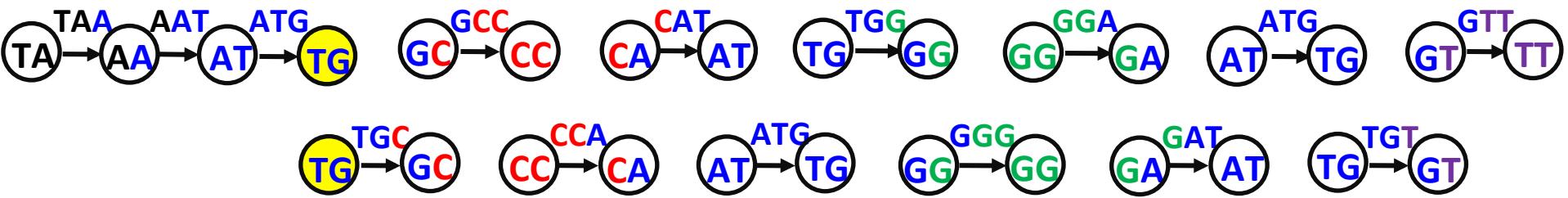
Gluing Identically Labeled Nodes

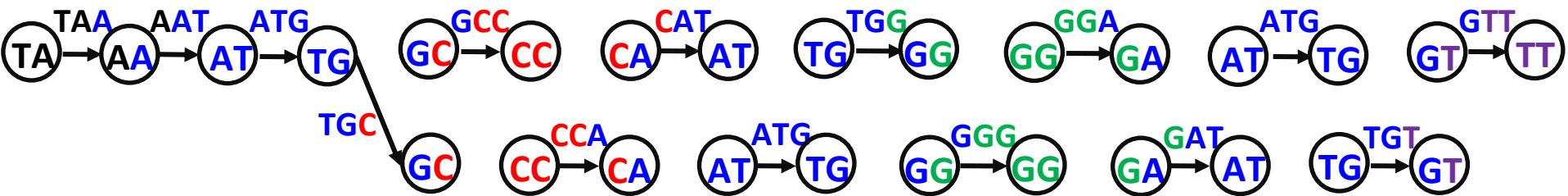


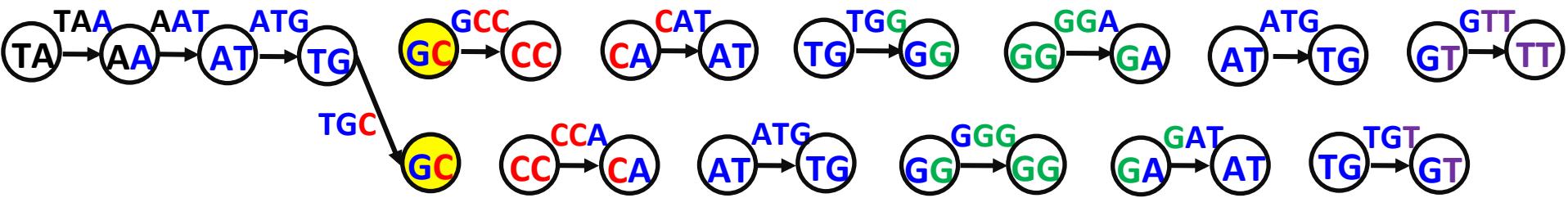


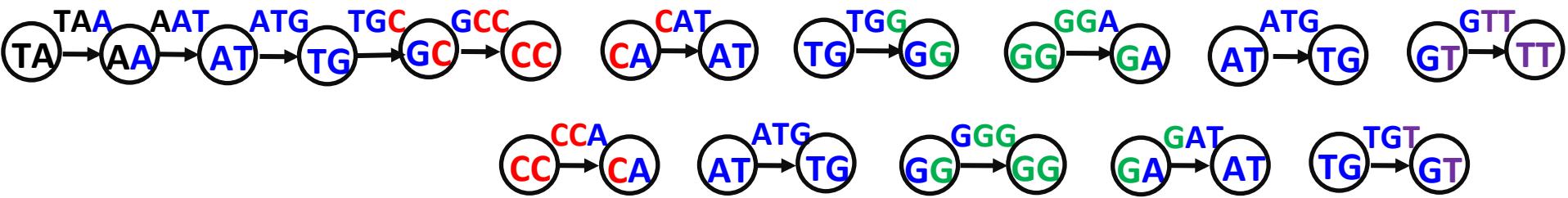


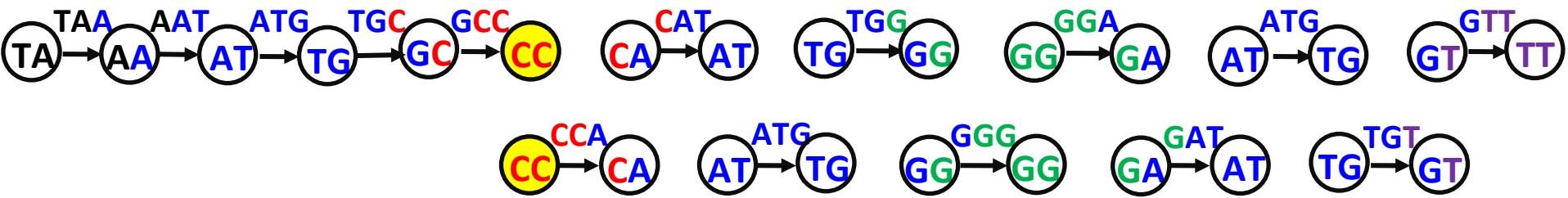


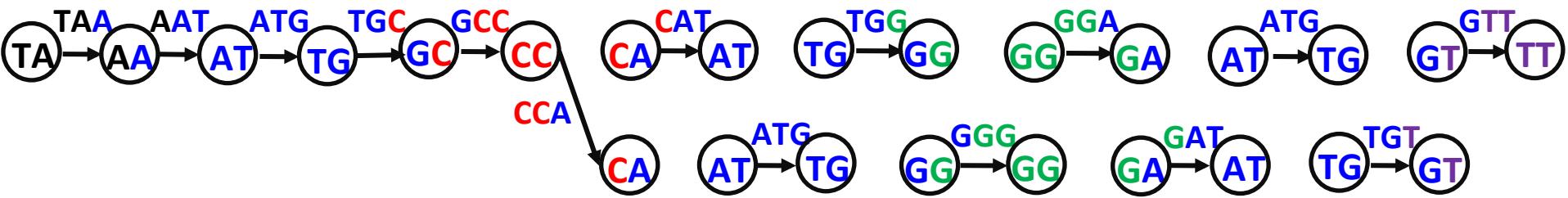


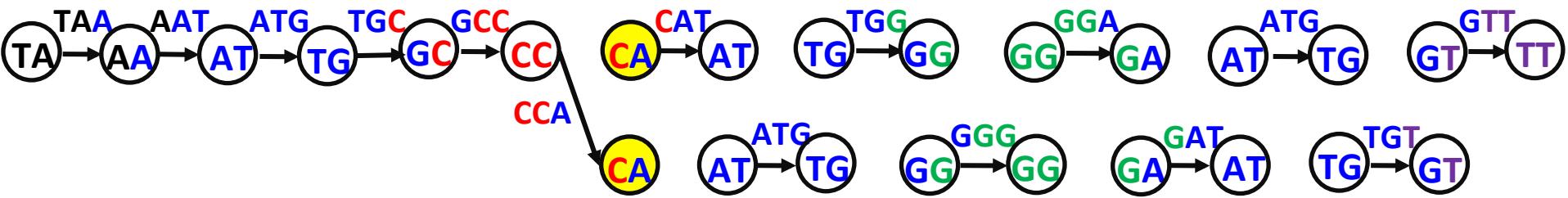


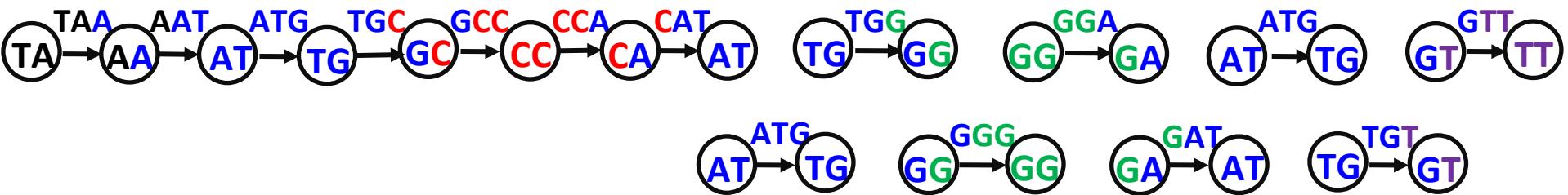


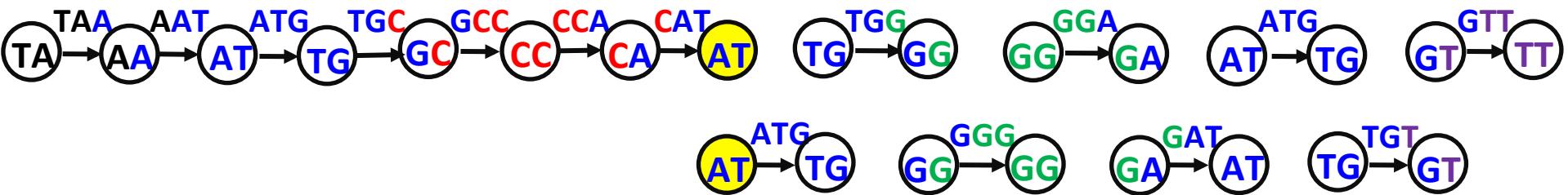


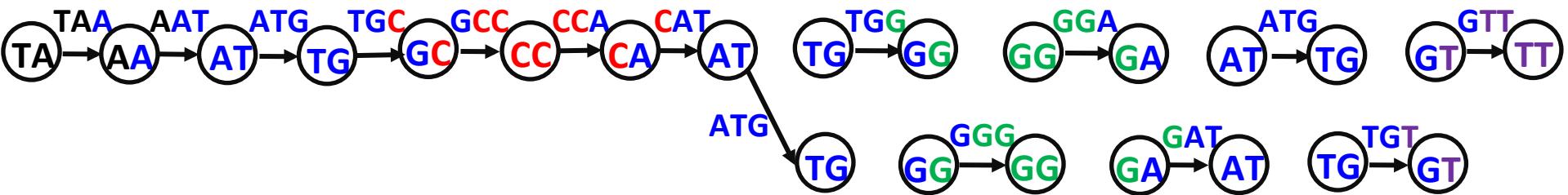


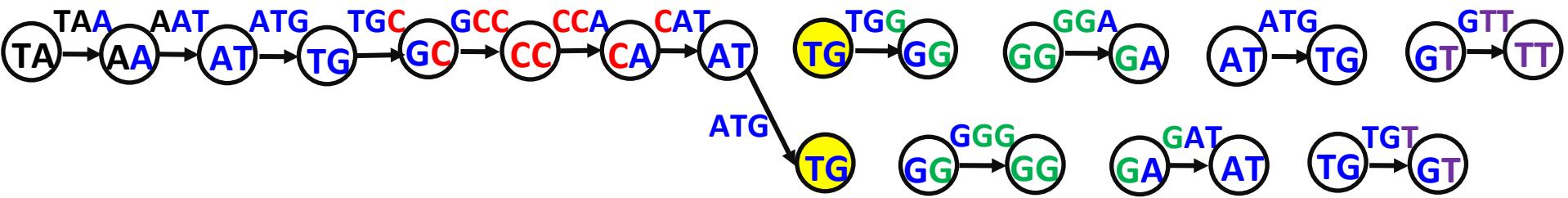


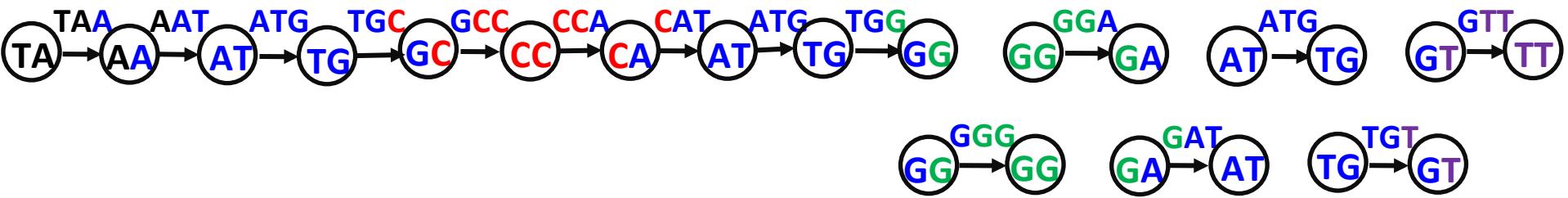


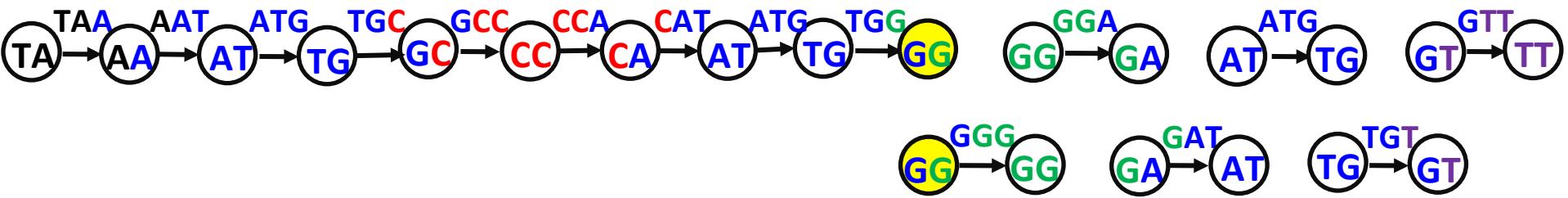


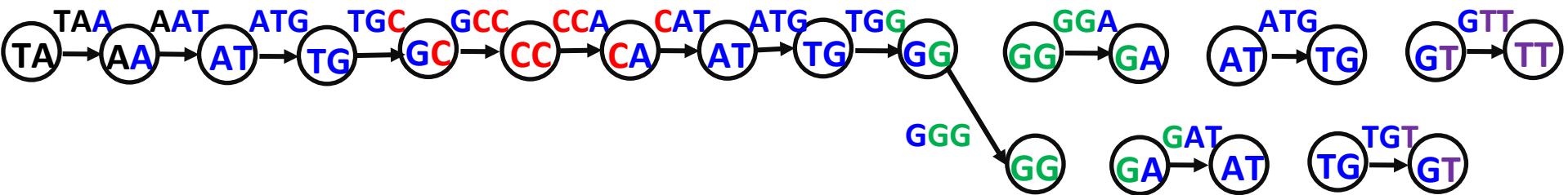


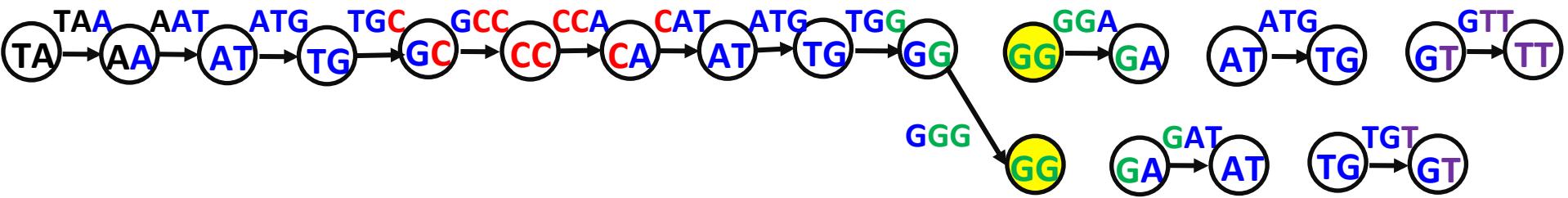


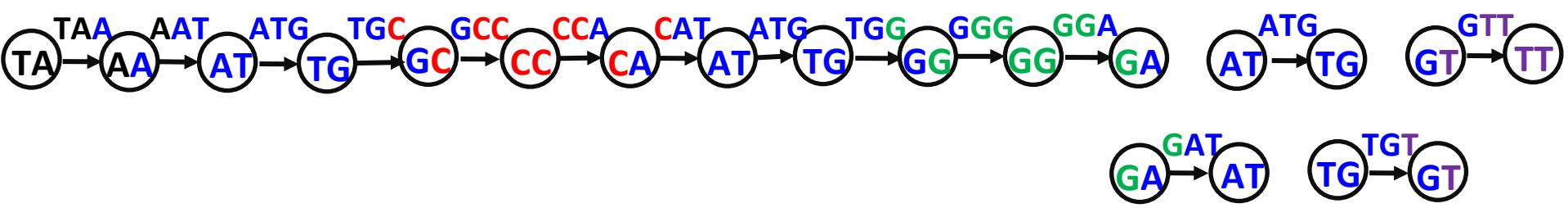


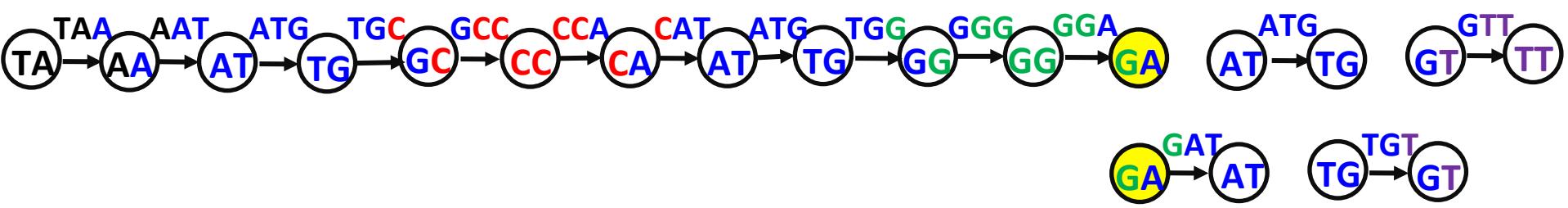


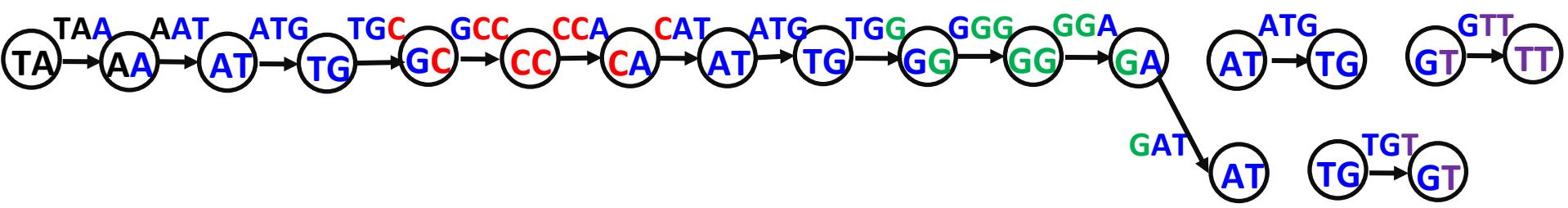


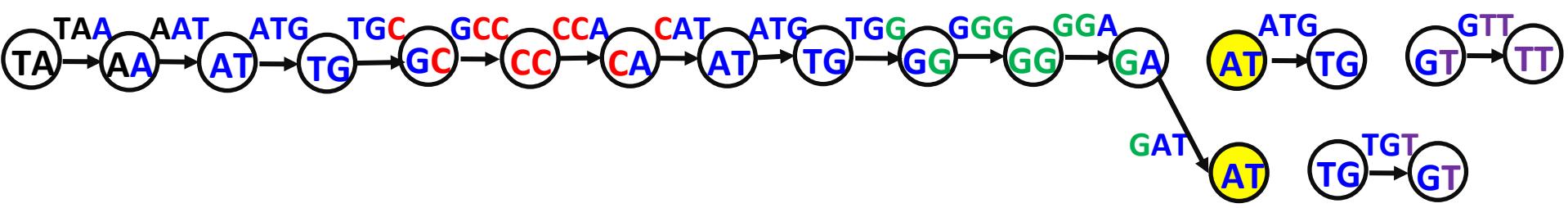


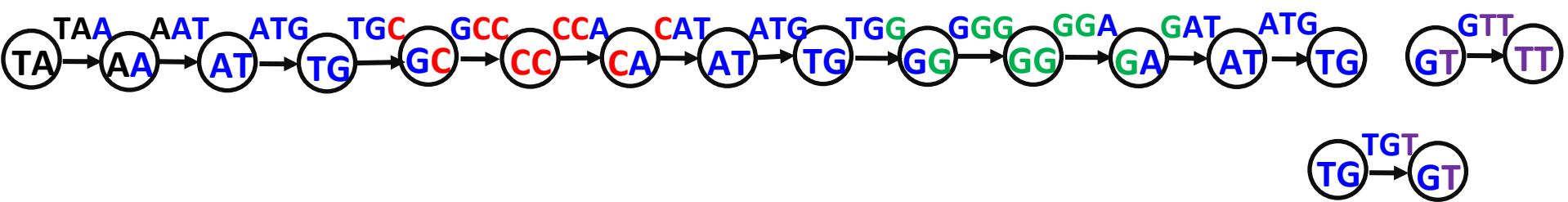


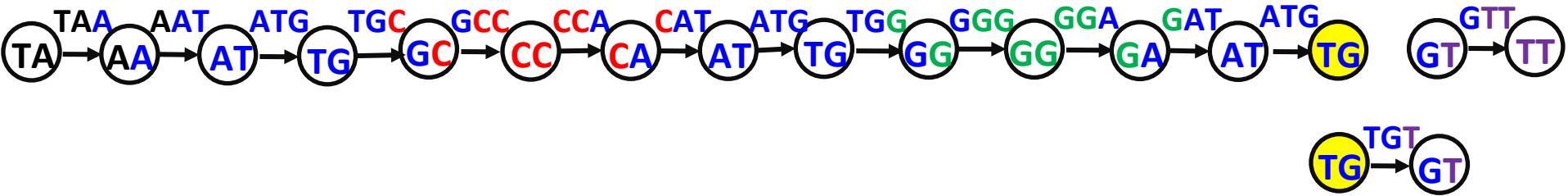


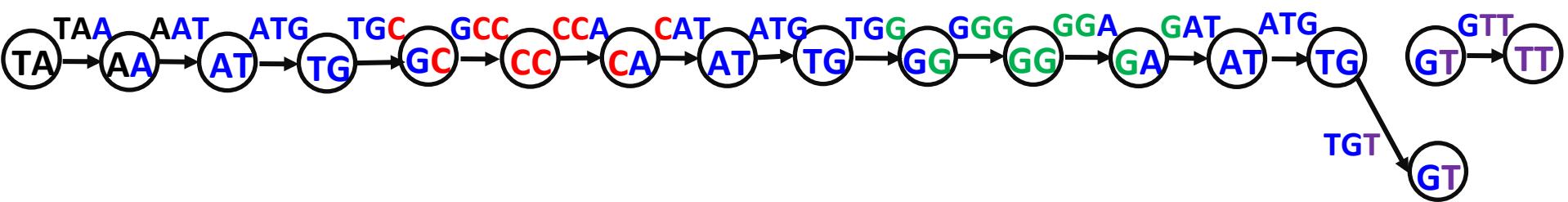


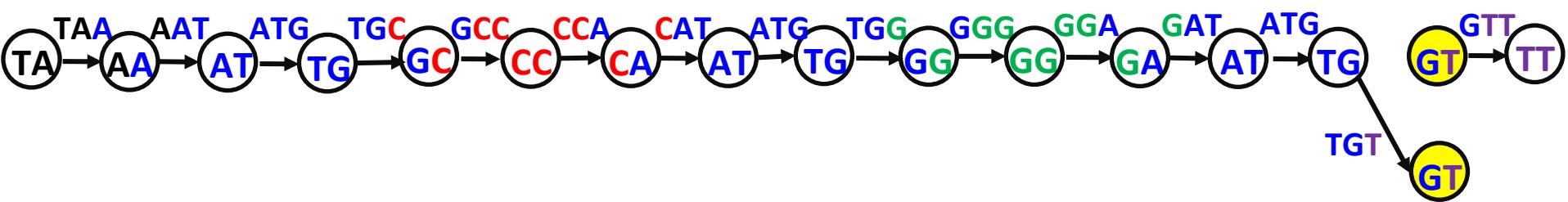








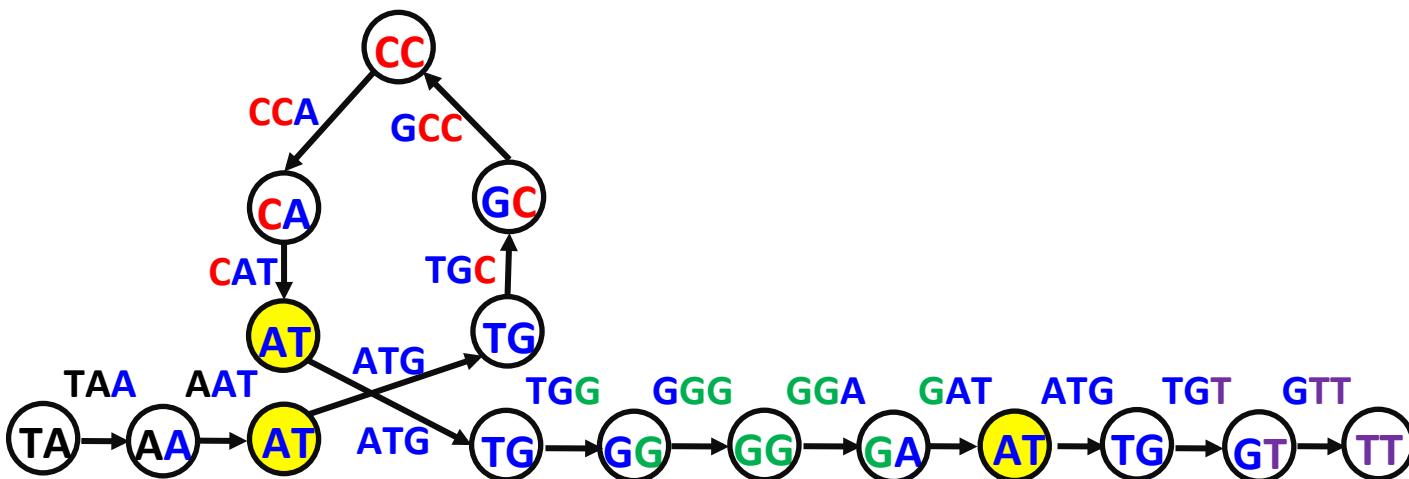
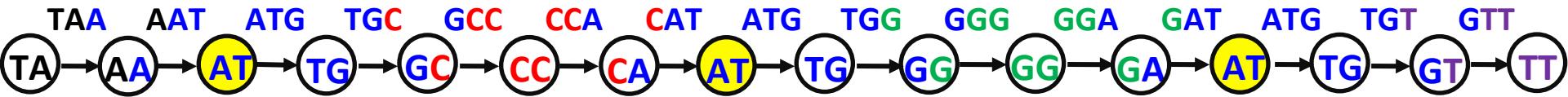




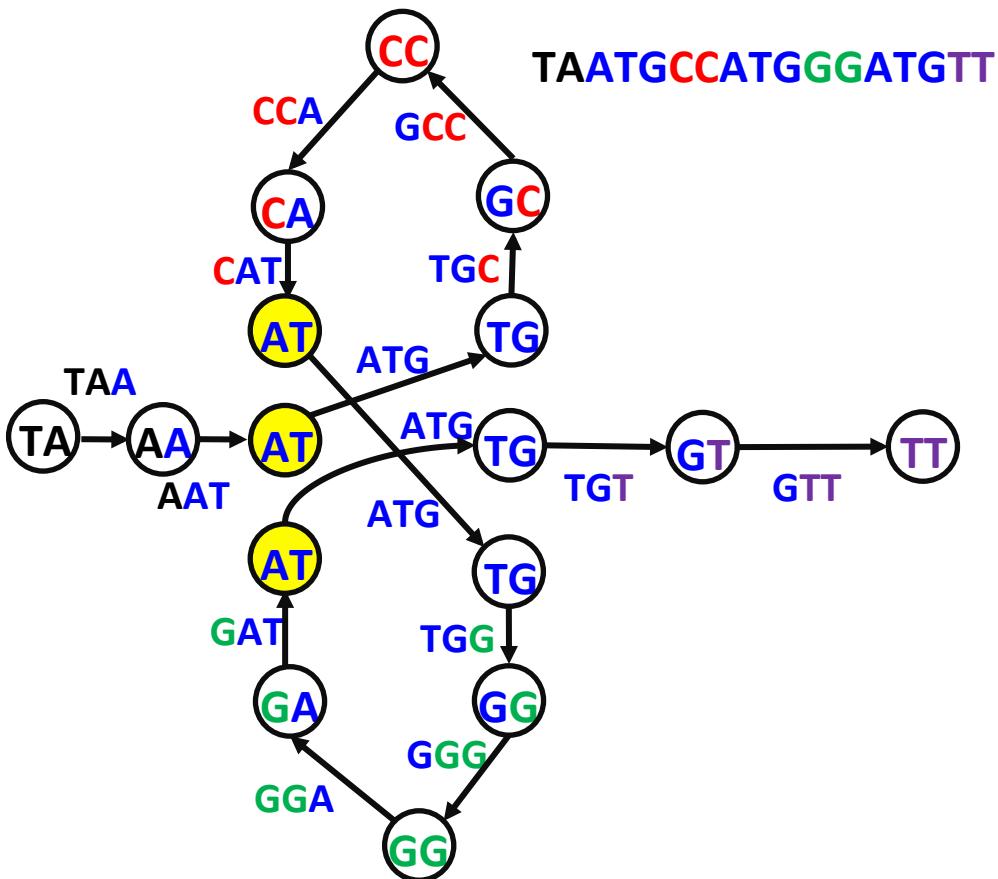
We Are Not Done with Gluing Yet

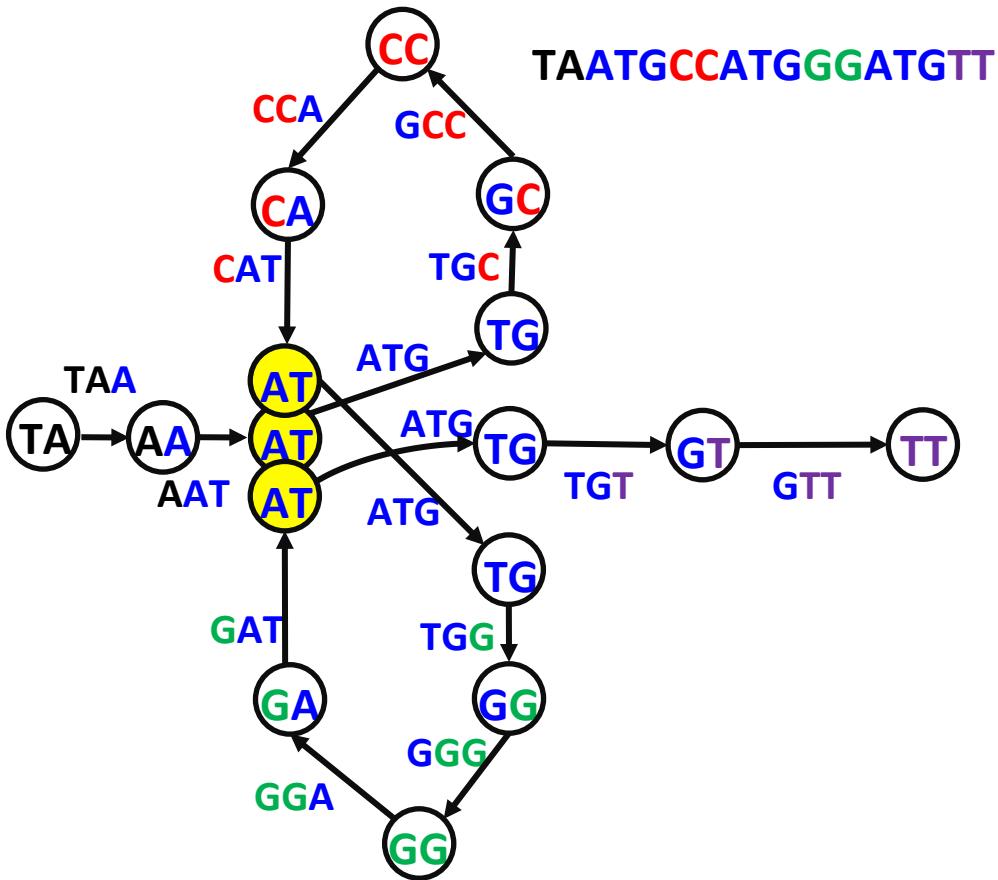


Gluing Identically Labeled Nodes

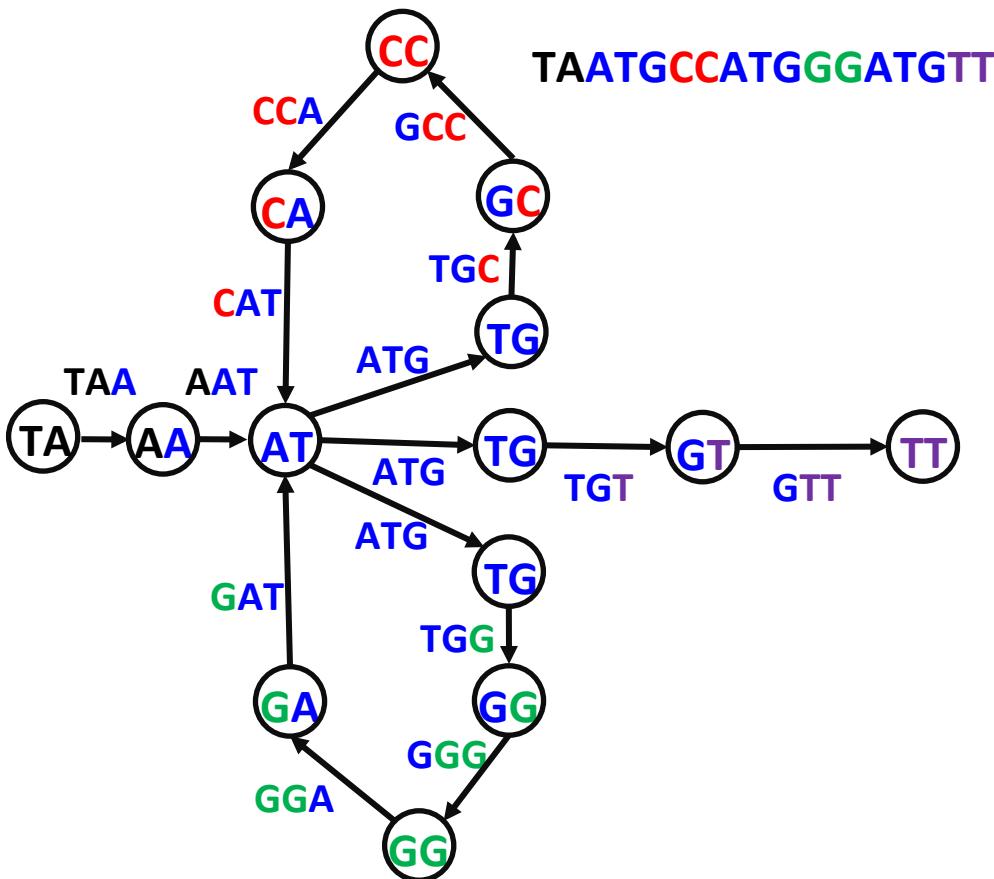


Gluing Identically Labeled Nodes

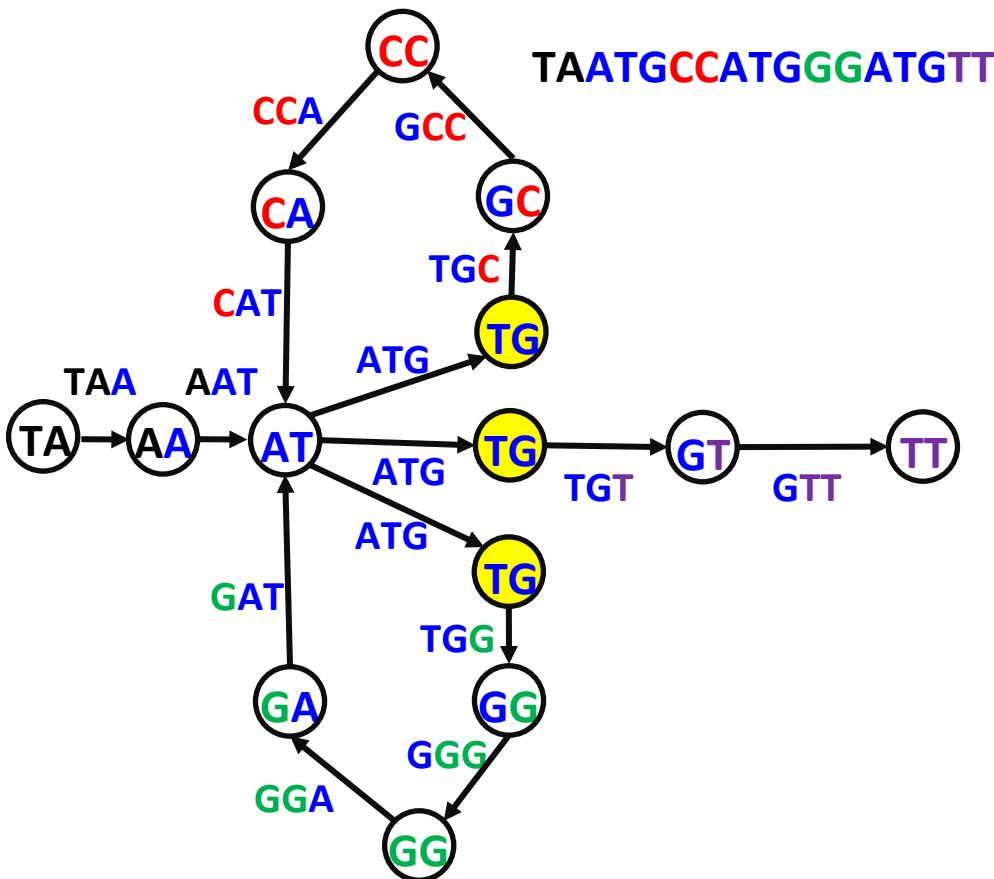




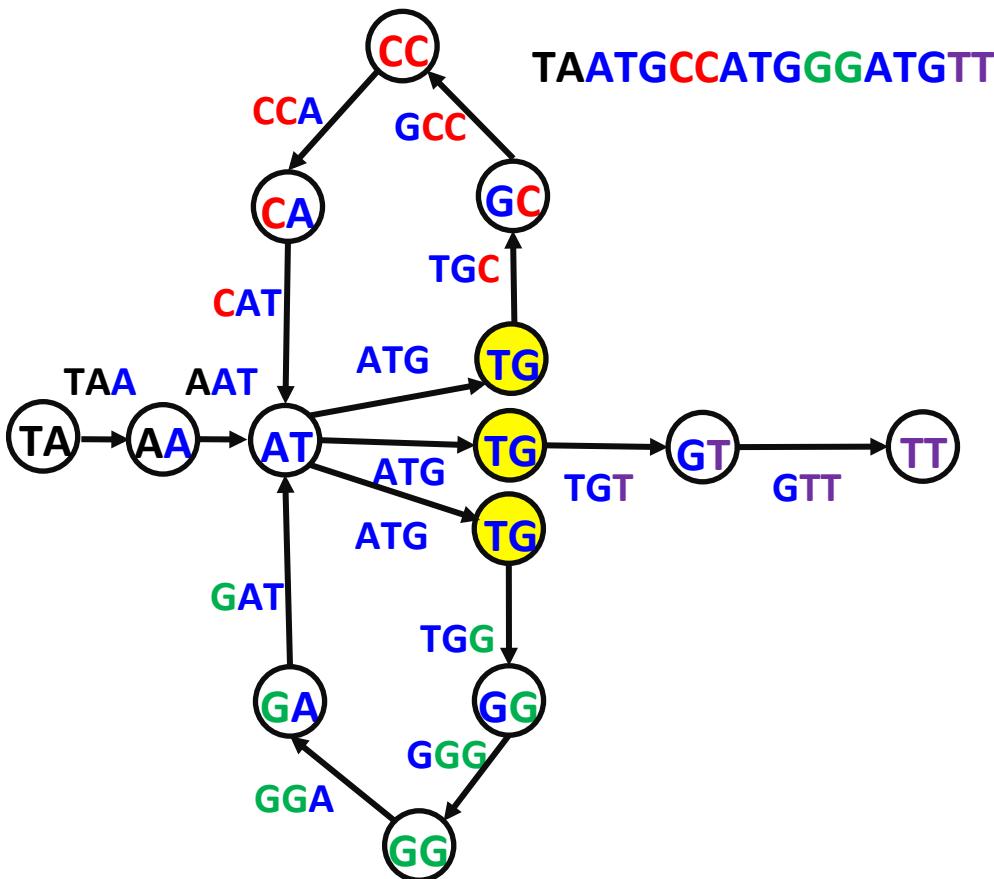
Gluing Identically Labeled Nodes



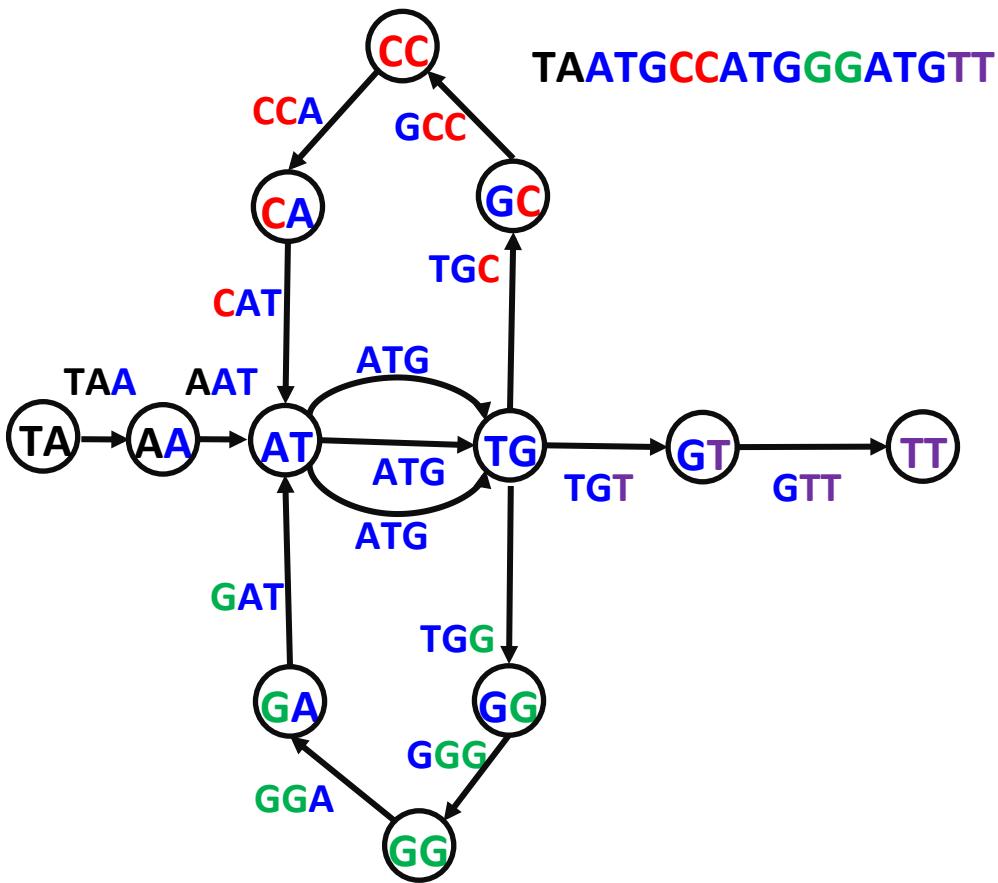
Gluing Identically Labeled Nodes



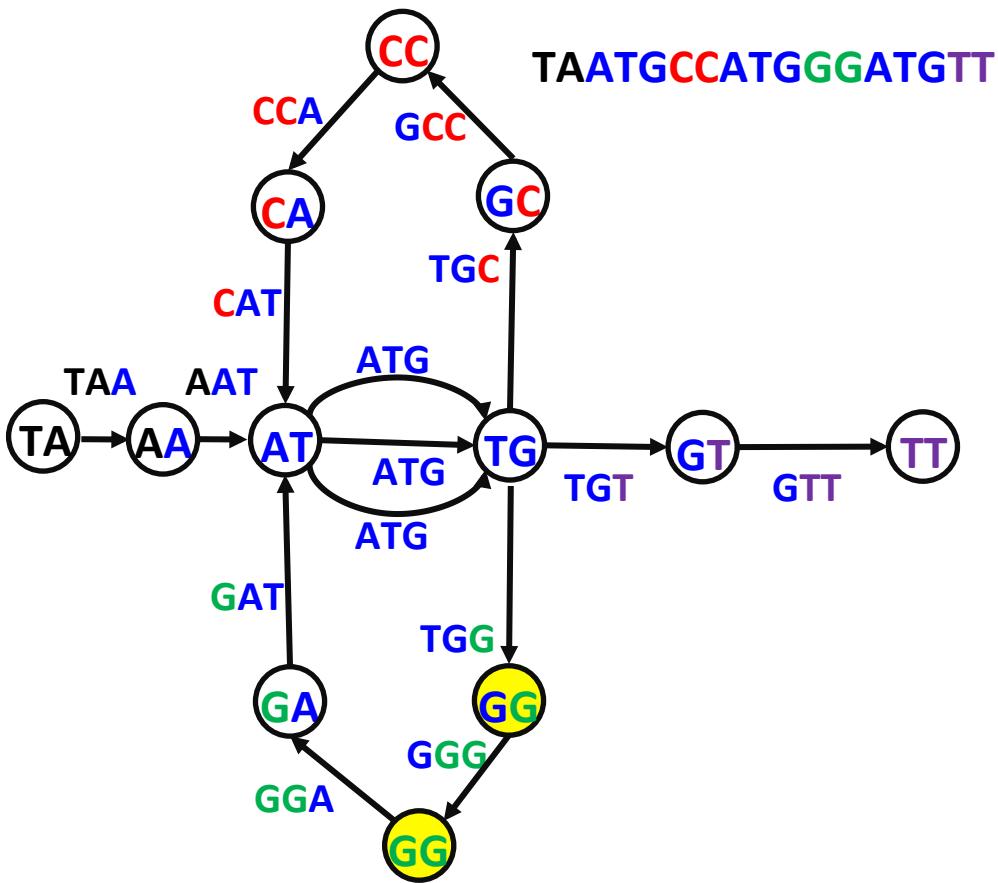
Gluing Identically Labeled Nodes



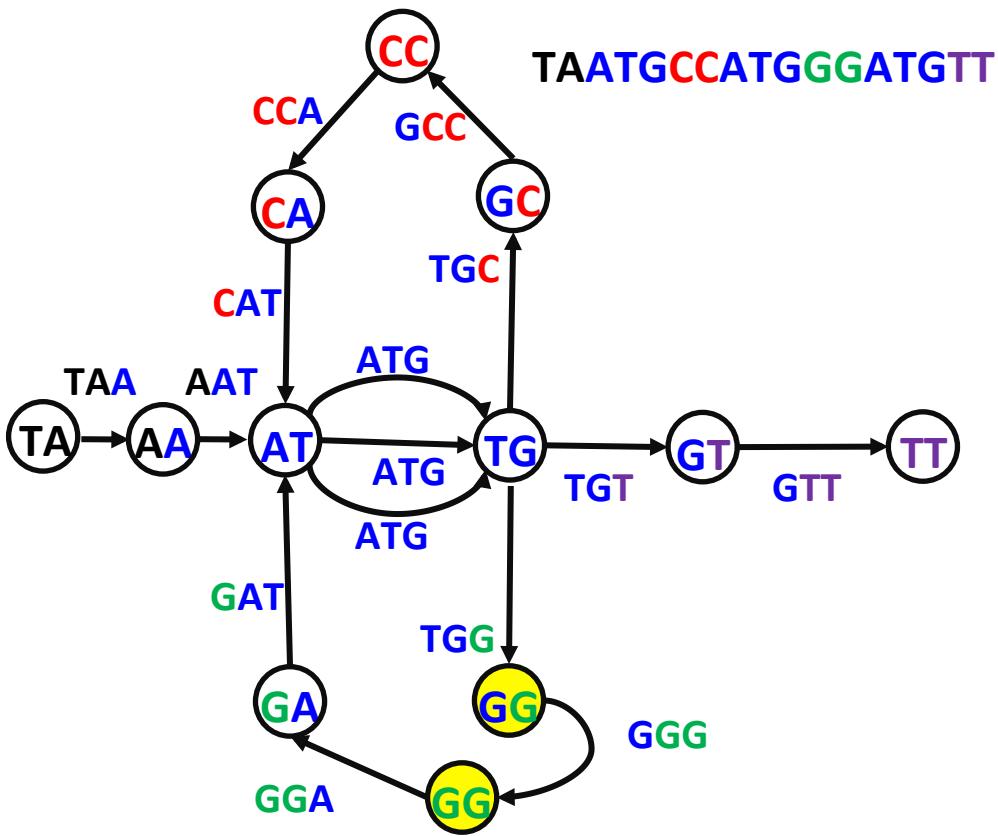
Gluing Identically Labeled Nodes



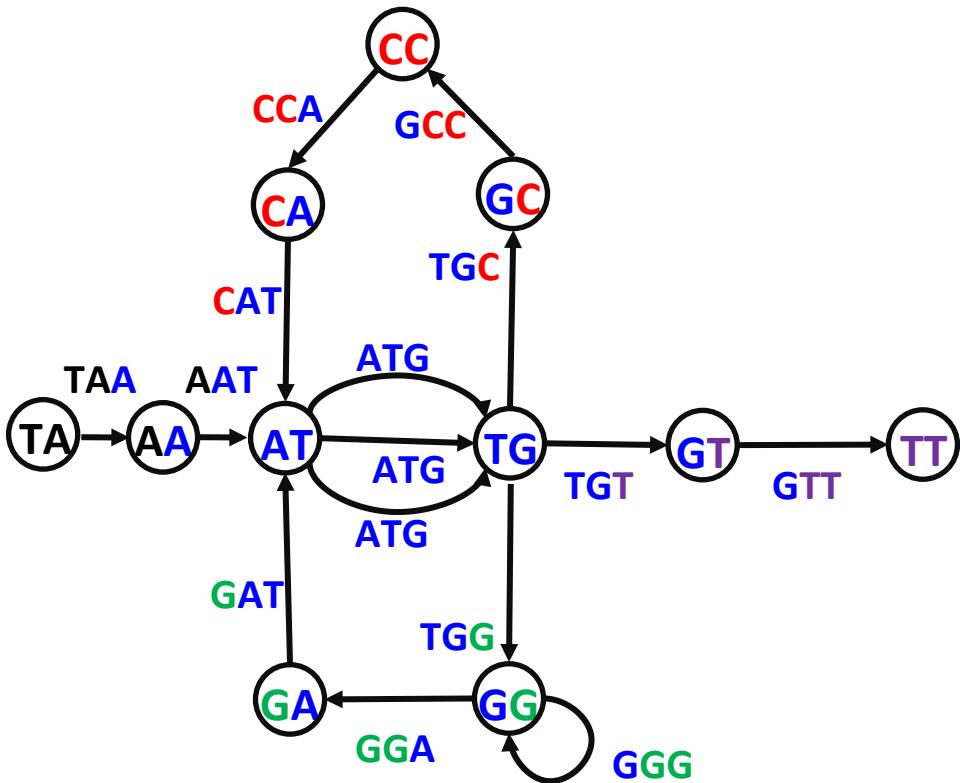
Gluing Identically Labeled Nodes



Gluing Identically Labeled Nodes



The Same de Bruijn Graph: DeBruin(Genome)=DeBruin(Genome Composition)



Constructing de Bruijn Graph

De Bruijn graph of a collection of k -mers:

- Represent every k -mer as an edge between its prefix and suffix
- Glue **ALL** nodes with identical labels.

DeBruijn(k -mers)

form a node for each $(k-1)$ -mer from k -mers

for each k -mer in k -mers

connect its prefix node with its suffix node by an edge

From Hamilton to Euler to de Bruijn



Universal String Problem (De Bruijn, 1946). Find a circular string containing each binary k-mer exactly once.

From Hamilton



to Euler

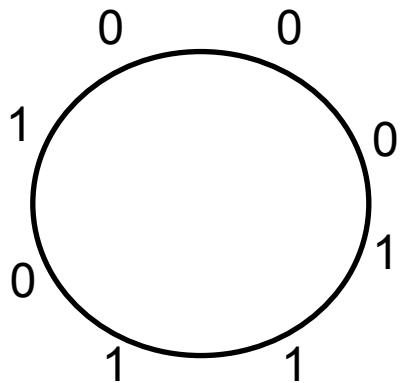


to de Bruijn



Universal String Problem (De Bruijn, 1946). Find a circular string containing each binary k-mer exactly once.

000 001 010 011 100 101 110 111



From Hamilton



to Euler

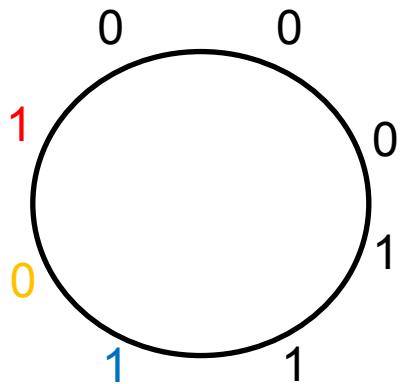


to de Bruijn



Universal String Problem (Nicolaas de Bruijn, 1946). Find a circular string containing each binary k-mer exactly once.

000 001 010 011 100 101 110 111



From Hamilton



to Euler

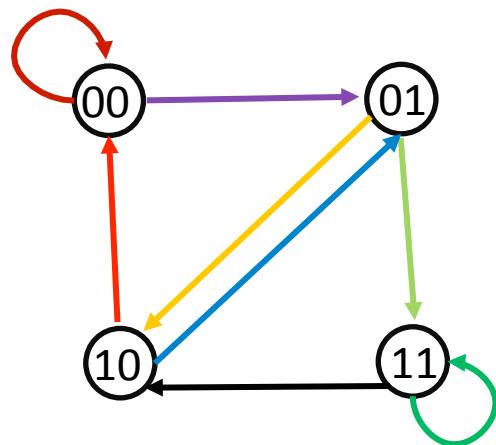
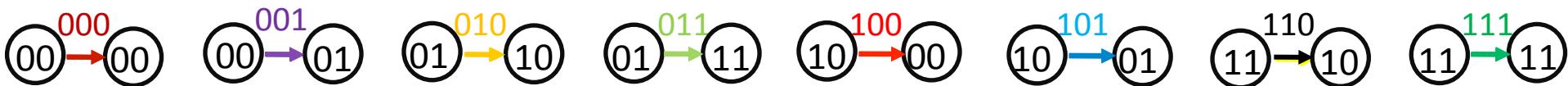


to de Bruijn



Universal String Problem (Nicolaas de Bruijn, 1946). Find a circular string containing each binary k-mer exactly once.

000 001 010 011 100 101 110 111



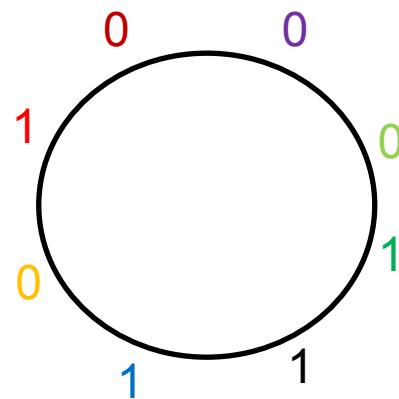
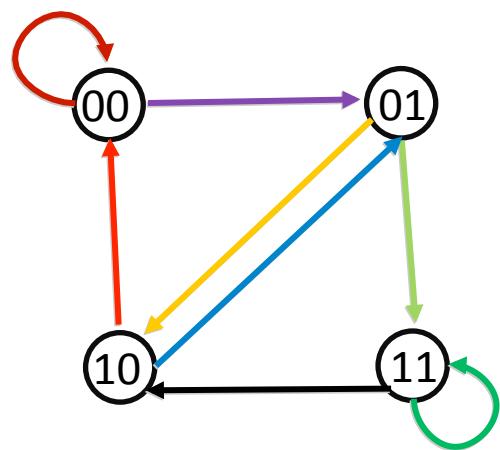
From Hamilton



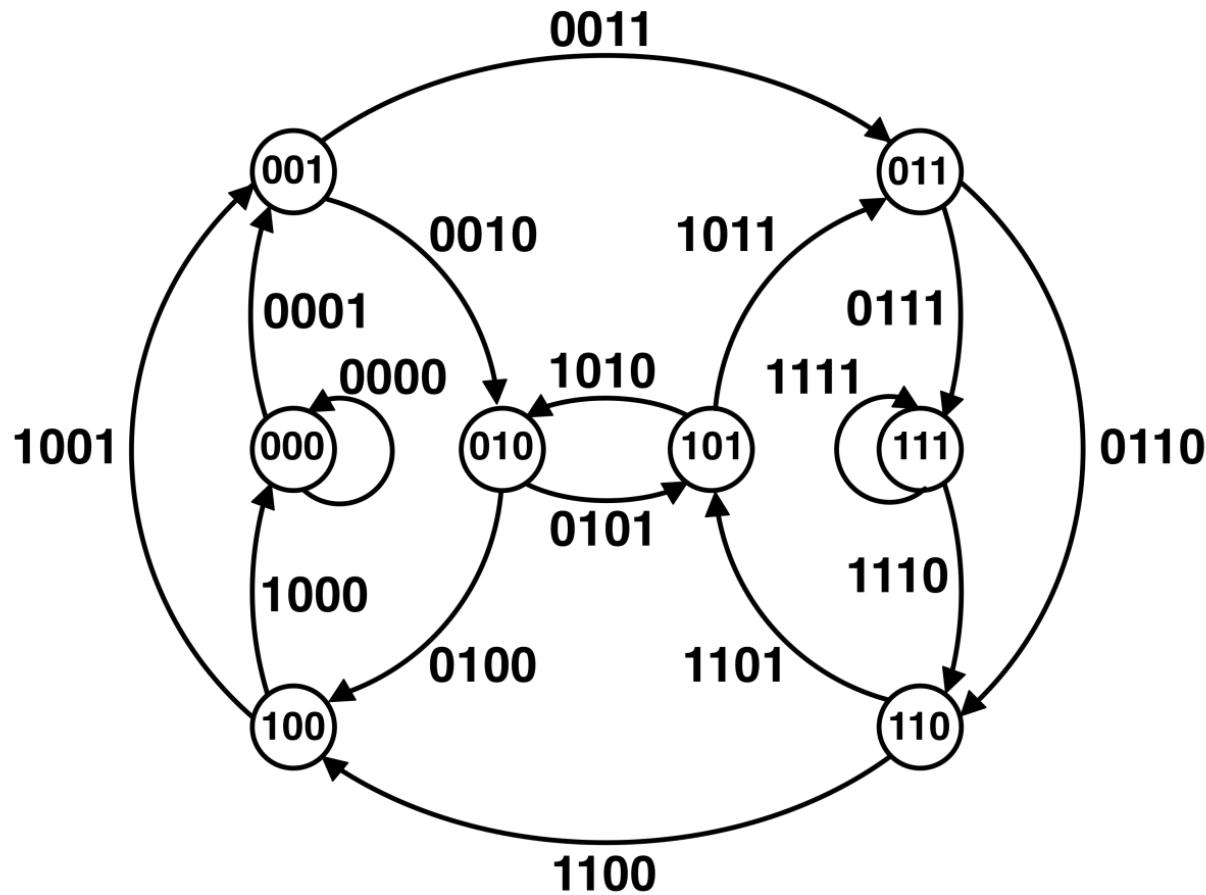
to Euler



to de Bruijn



De Bruijn Graph for 4-Universal String



Does it have an Eulerian cycle? If yes, how can we find it?

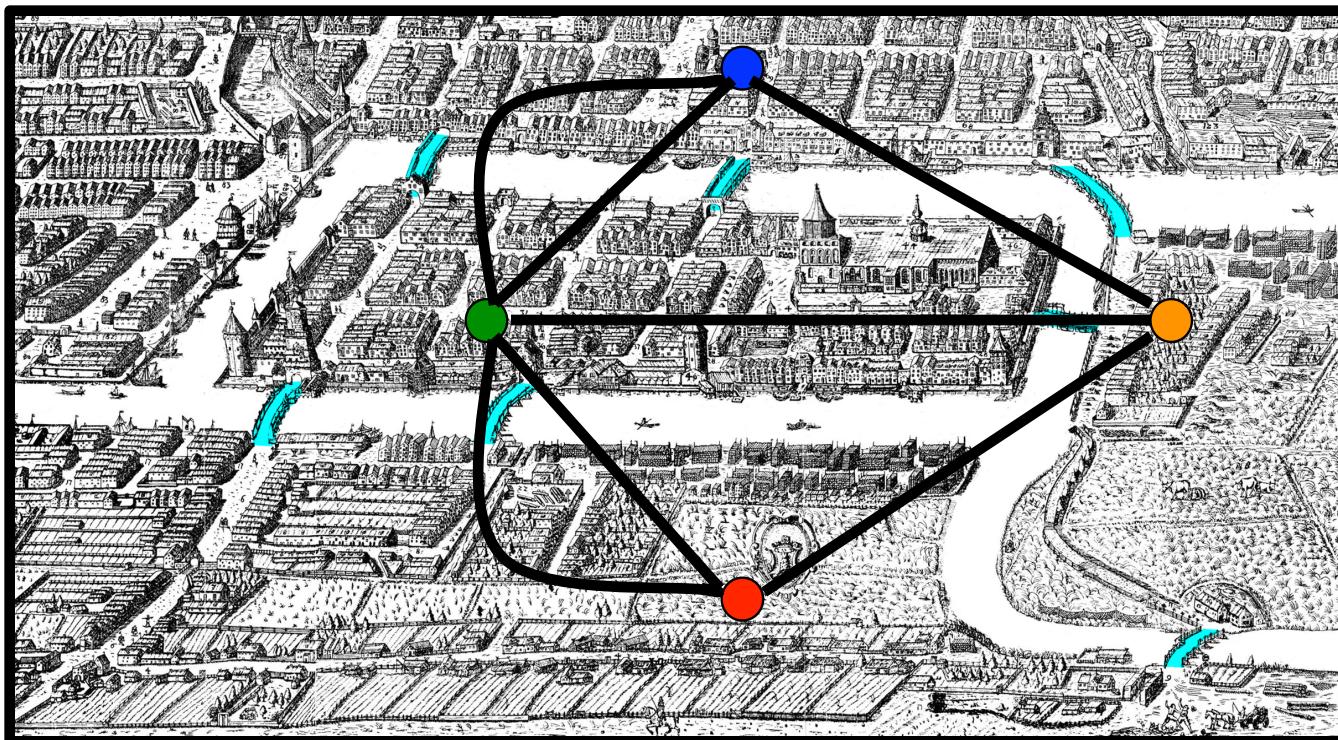
Outline

- What Is Genome Sequencing?
- Exploding Newspapers
- The String Reconstruction Problem
- String Reconstruction as a Hamiltonian Path Problem
- String Reconstruction as an Eulerian Path Problem
- Similar Problems with Different Fates
- De Bruijn Graphs
- Euler's Theorem
- Assembling Read-Pairs
- De Bruijn Graphs Face Harsh Realities of Assembly

Eulerian CYCLE Problem

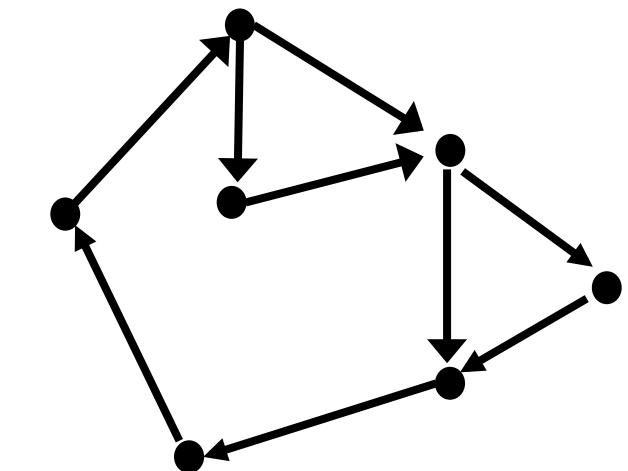
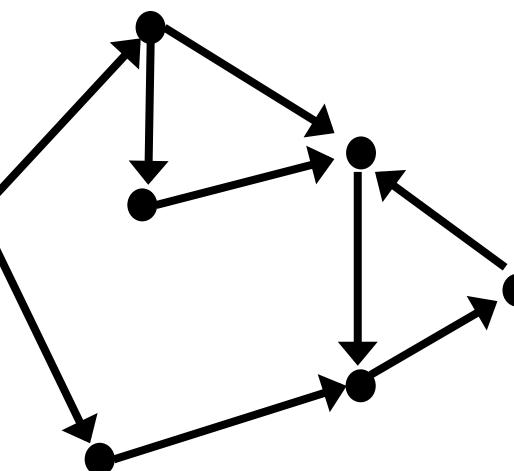
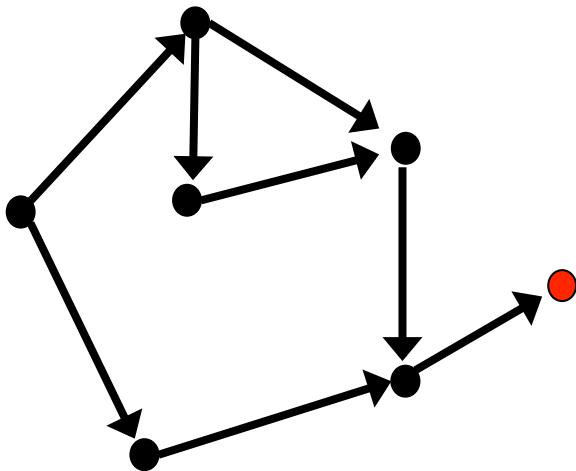
Eulerian CYCLE Problem. Find an Eulerian cycle in a graph.

- Input. A graph.
- Output. A cycle visiting every edge in the graph exactly once.



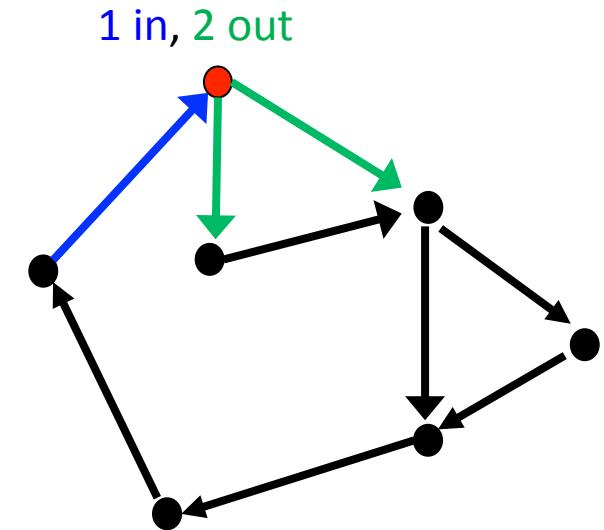
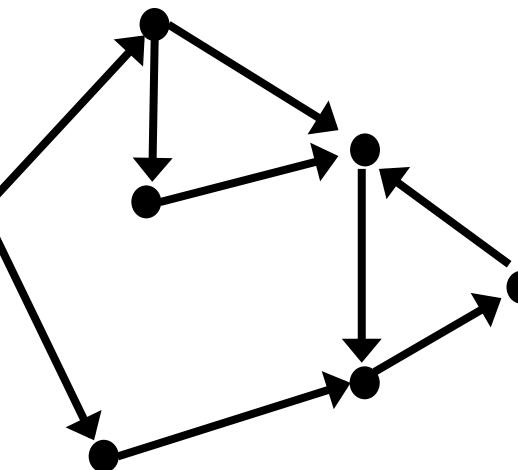
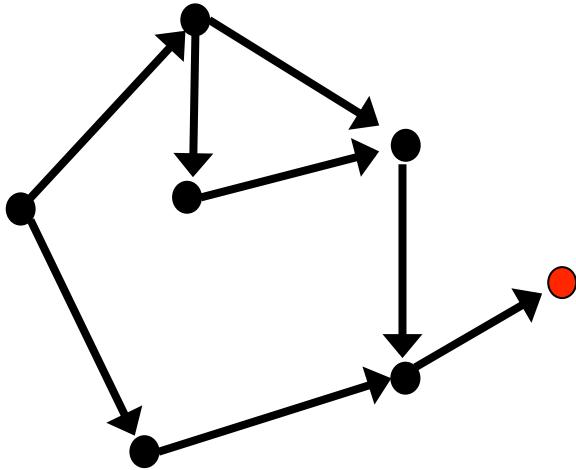
A Graph is Eulerian if It Contains an Eulerian Cycle.

Is this graph Eulerian?



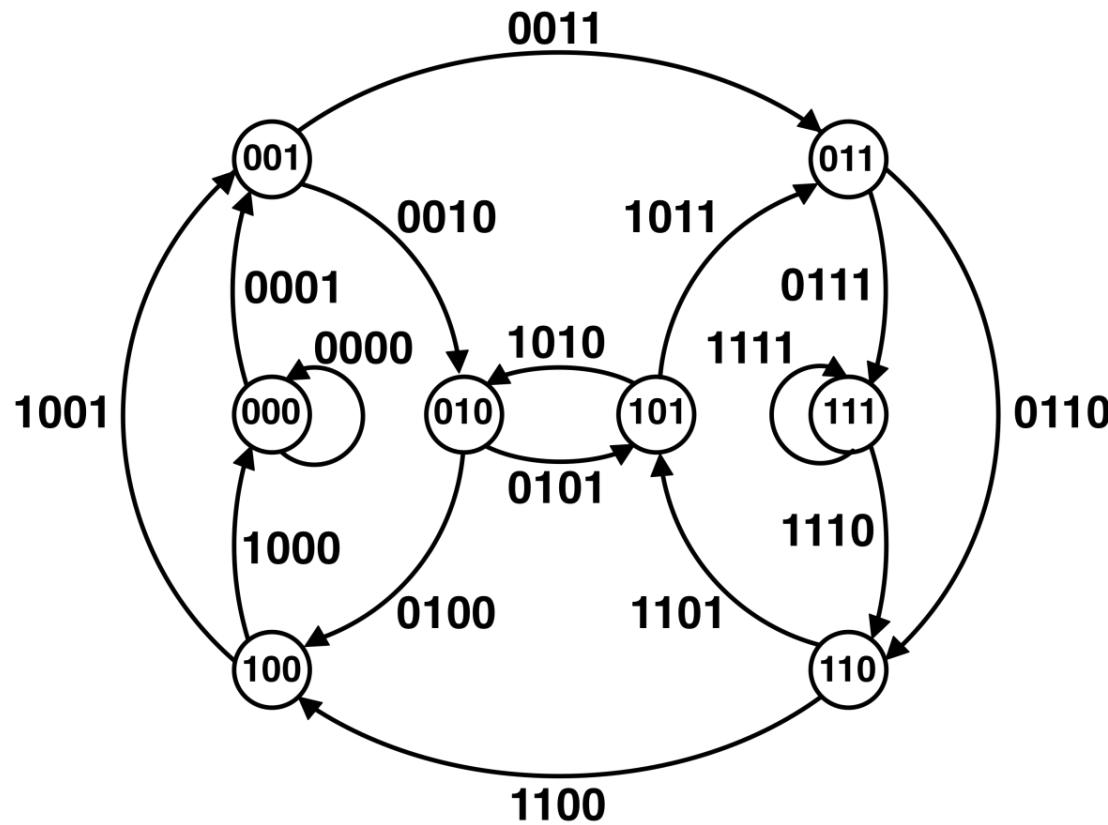
A Graph is Eulerian if It Contains an Eulerian Cycle.

Is this graph Eulerian?



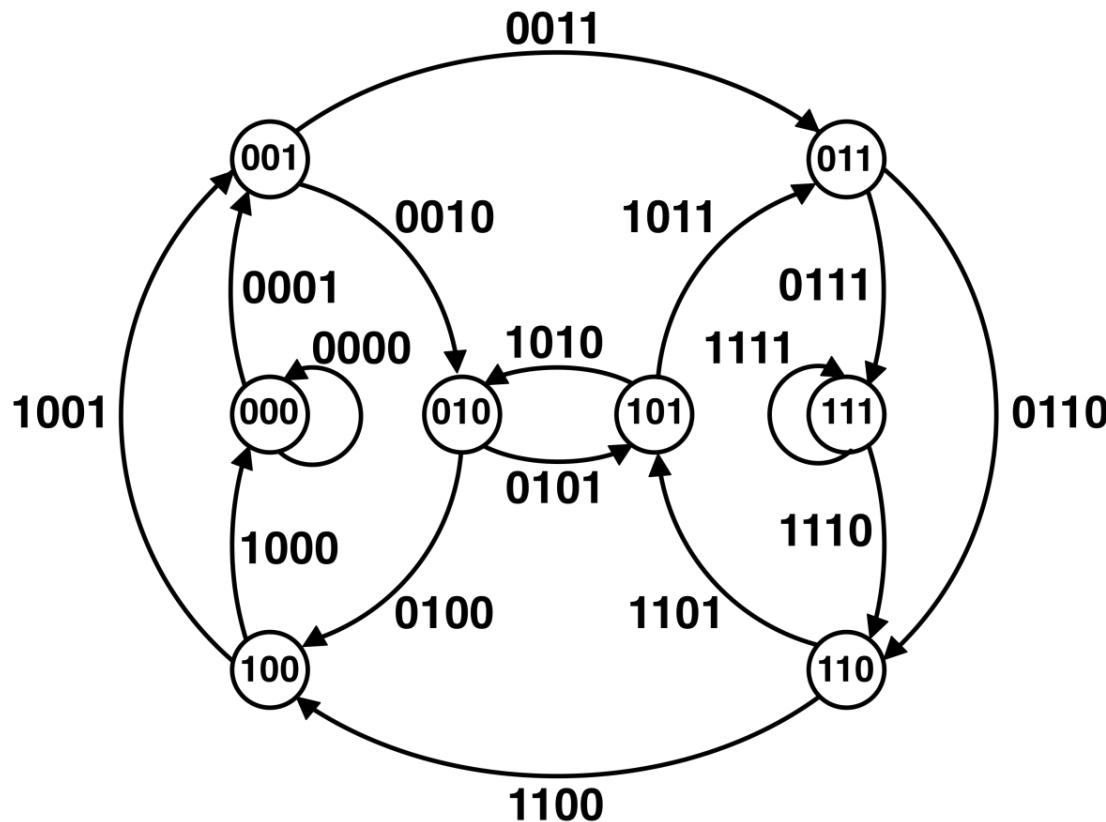
A graph is balanced if **indegree = outdegree** for each node

Is the Graph for 4-Universal String Balanced?



Euler's Theorem

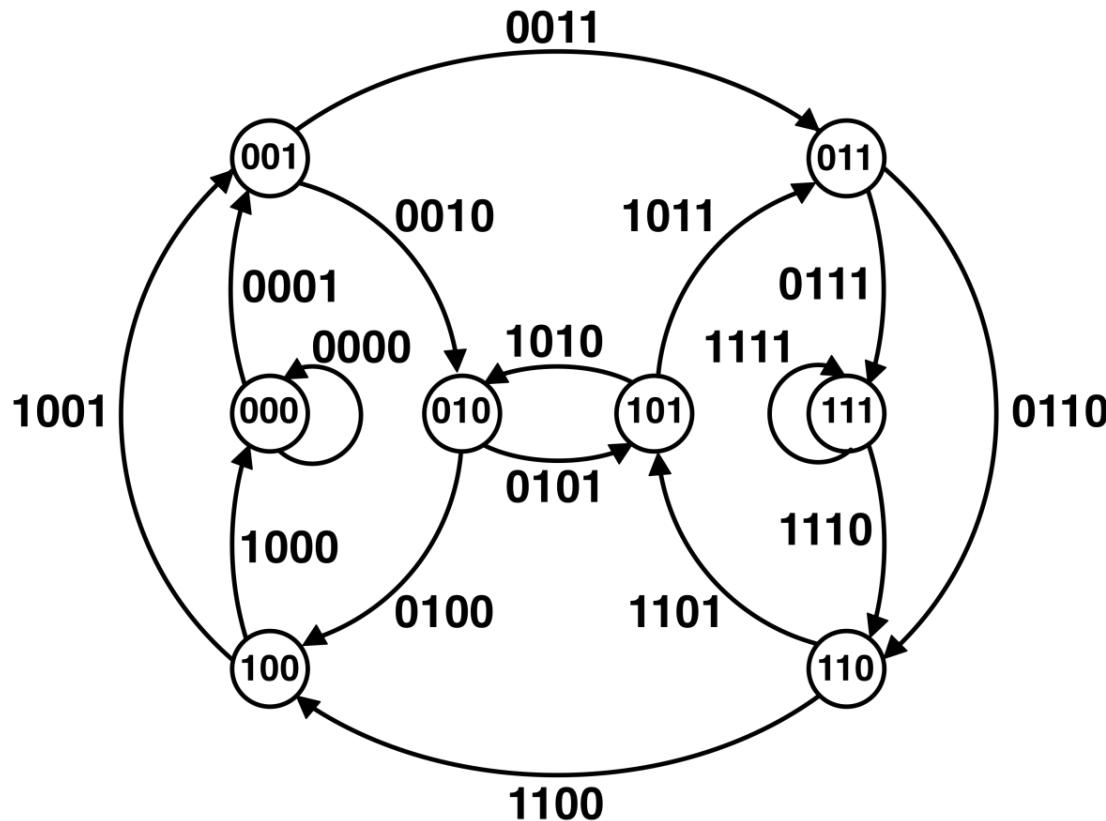
- Every Eulerian graph is balanced



Euler's Theorem



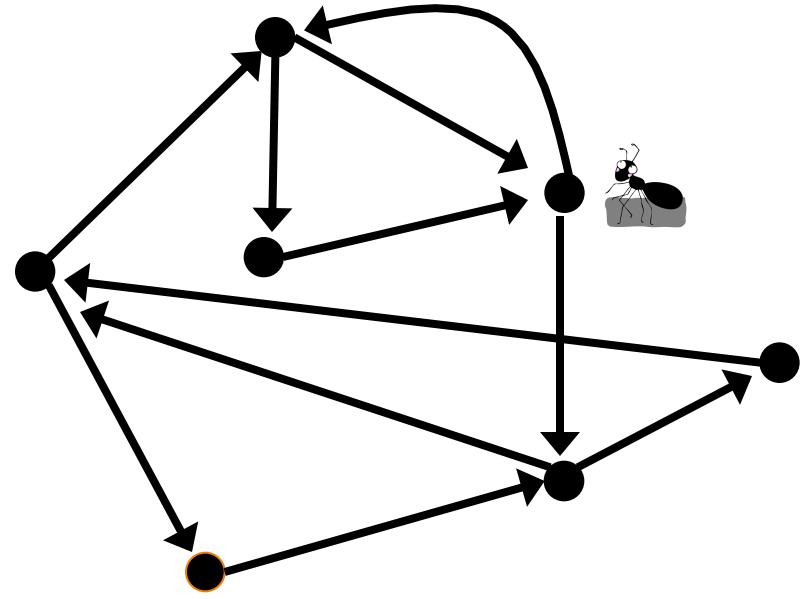
- Every Eulerian graph is balanced
 - **Every balanced* graph is Eulerian**



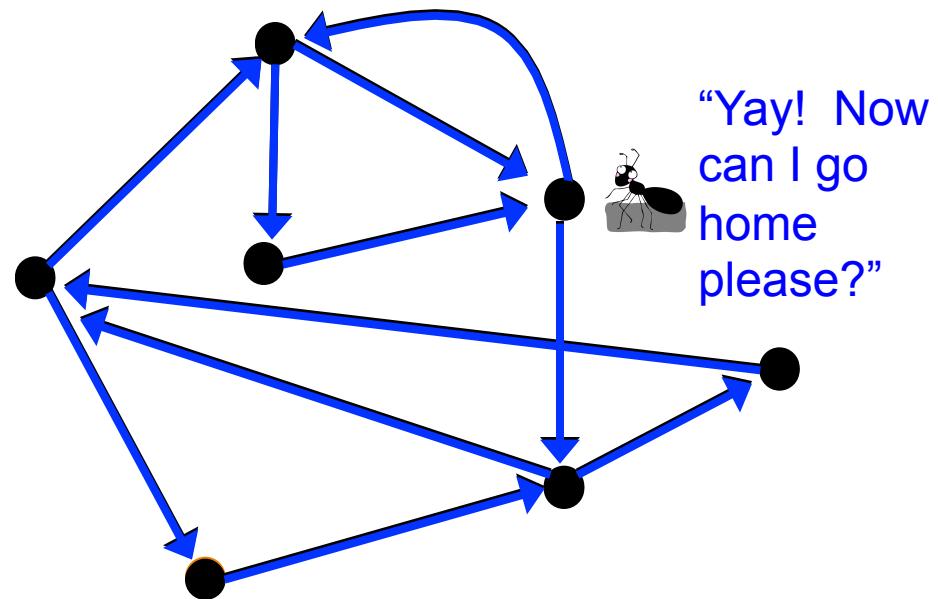
(*) and strongly connected, of course!

Recruiting an Ant to Prove Euler's Theorem

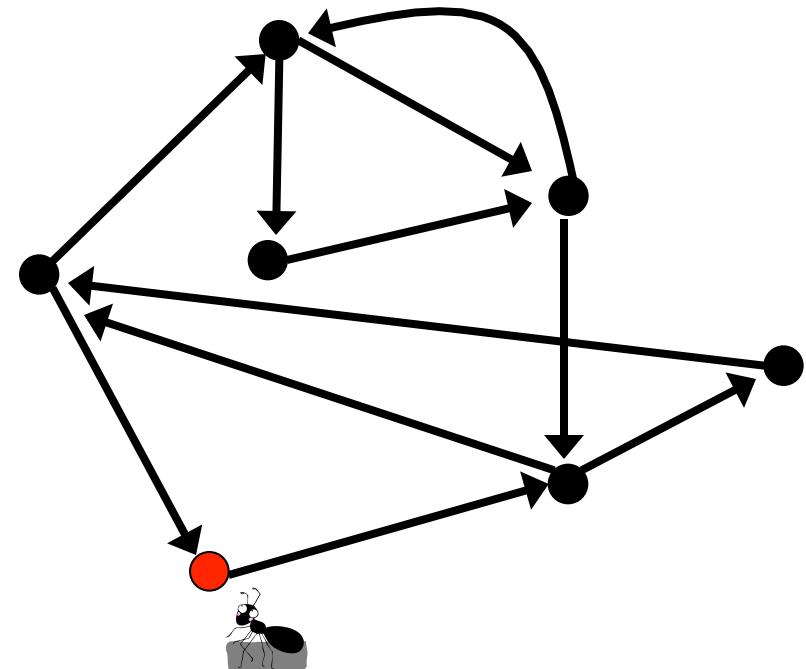
Let an ant randomly walk through the graph.
The ant cannot use the same edge twice!



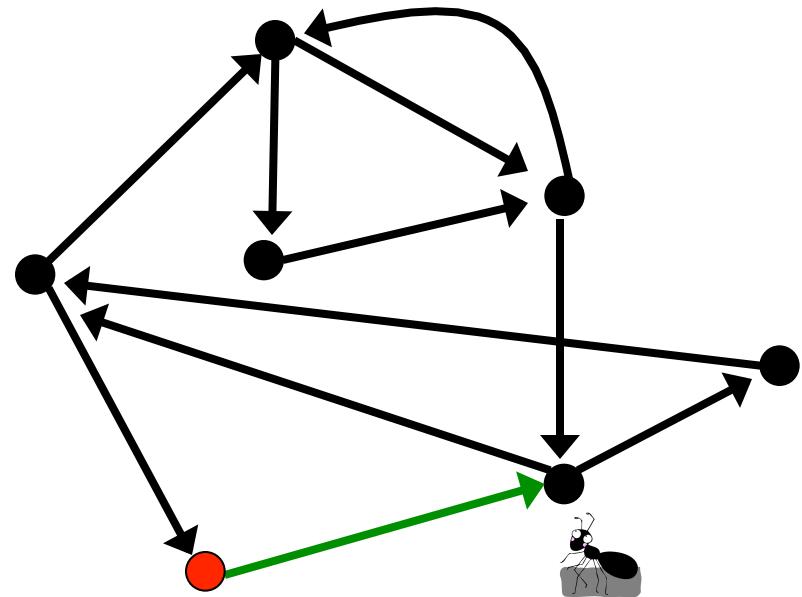
If Ant Was a Genius...



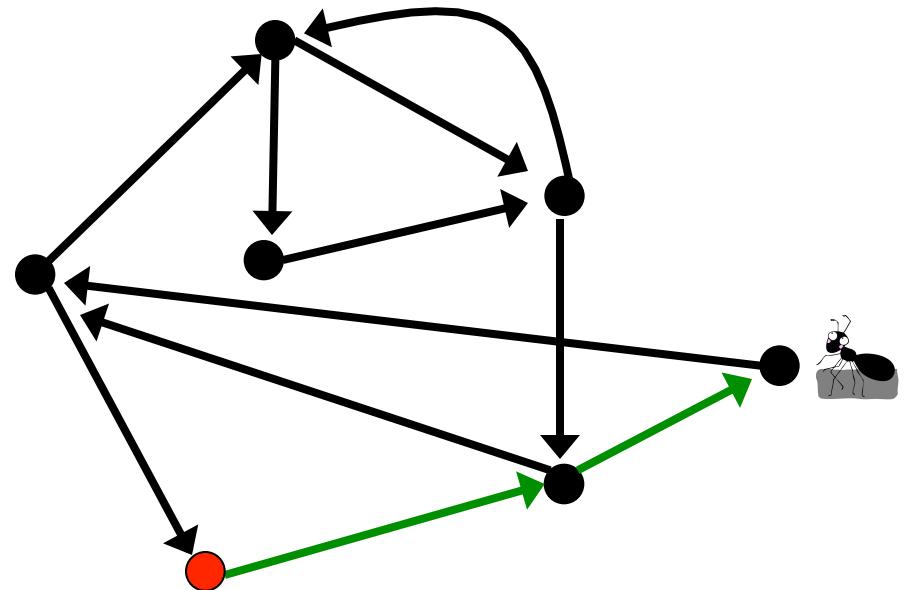
A Less Intelligent Ant Would Randomly Choose a Node and Start Walking...



Walking...

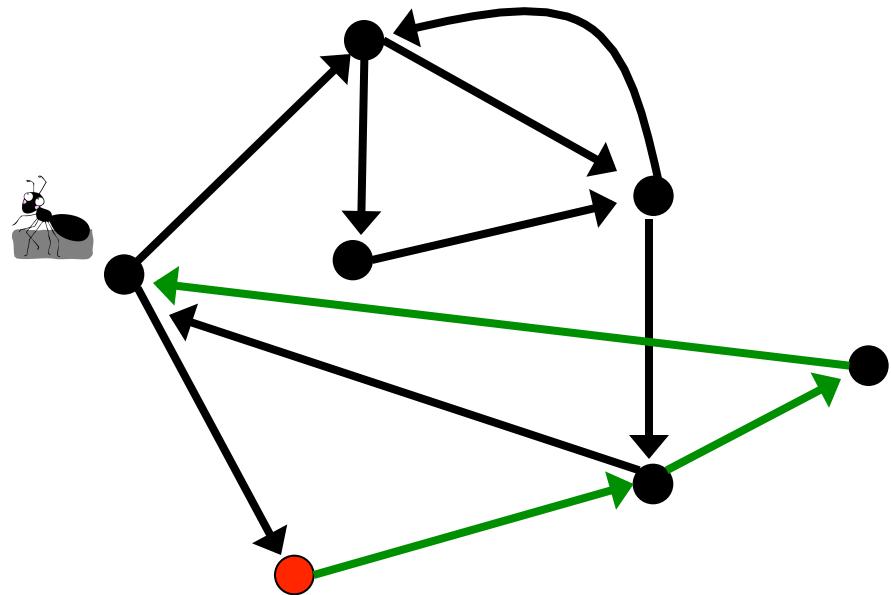


Walking... and Walking...

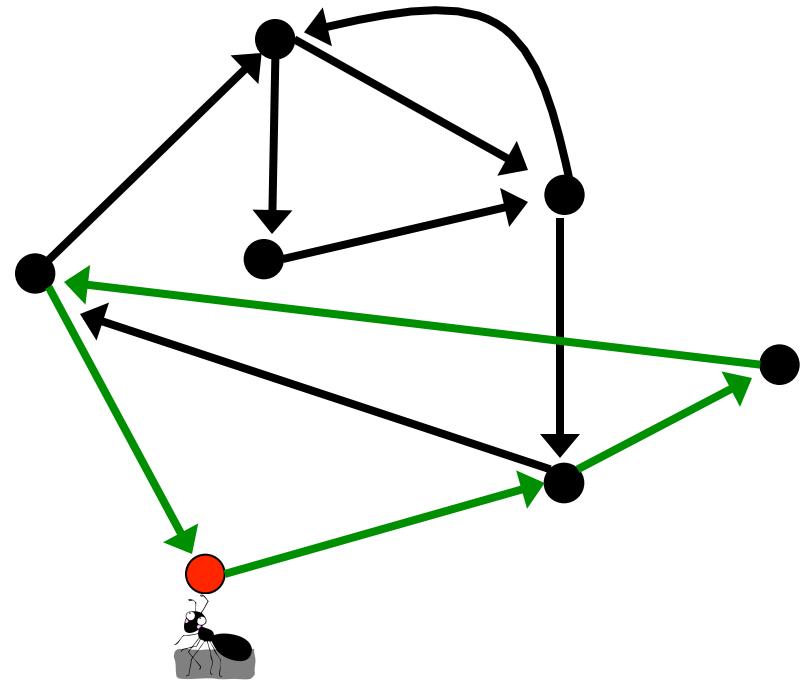


Walking... and Walking... and Walking...

Can it get stuck? In what node?

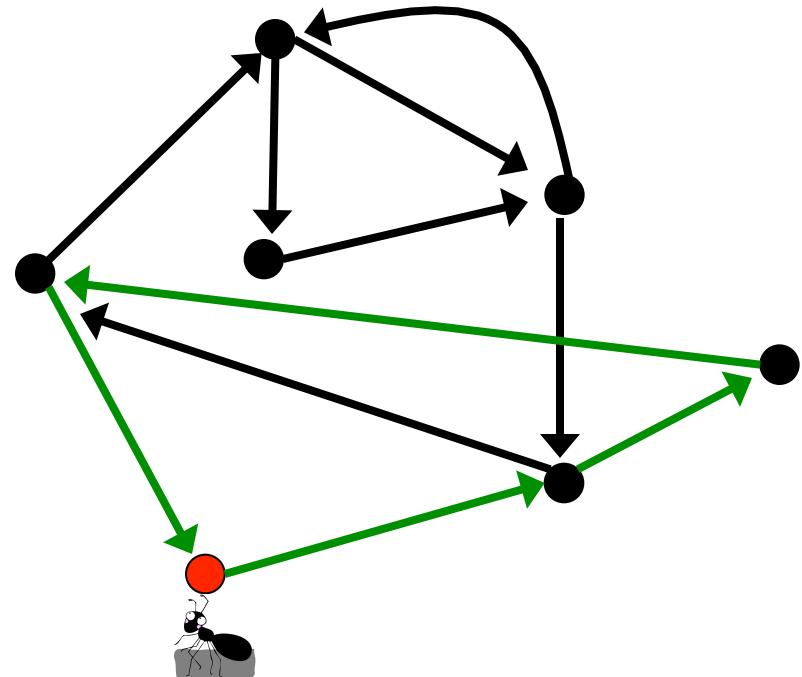


The Ant Can Only Get Stuck at the Starting Node



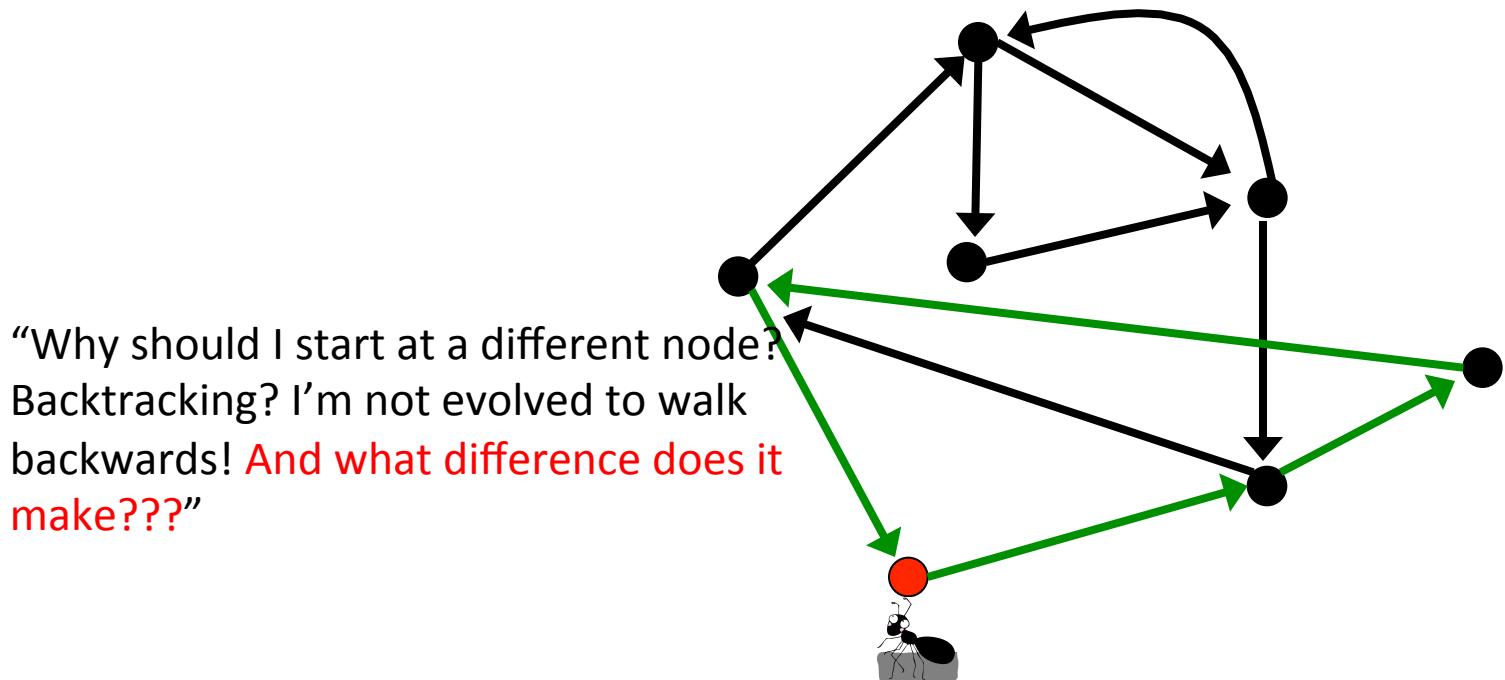
The Ant Has Completed a Cycle BUT has not Proven Euler's theorem yet...

The constructed cycle is not Eulerian. **Can we enlarge it?**



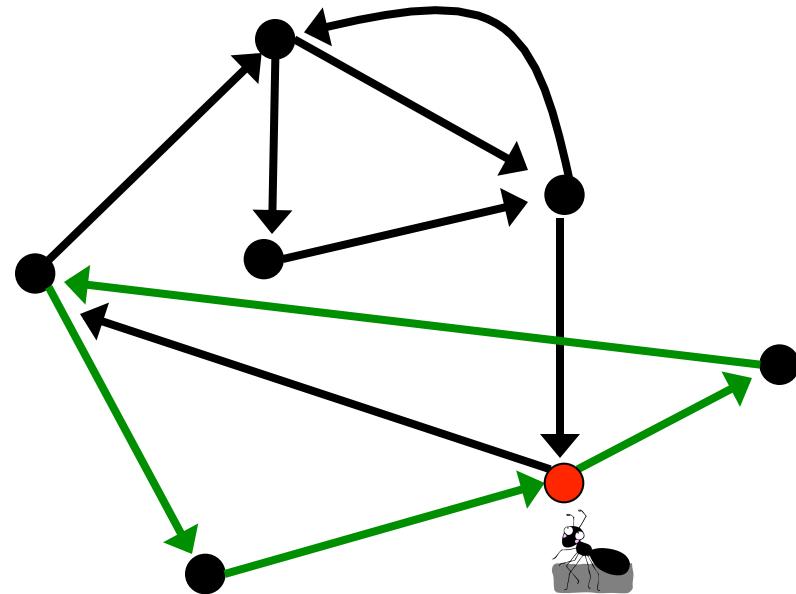
Let's Start at a Different Node in the Green Cycle

Let's start at a node with still unexplored edges.



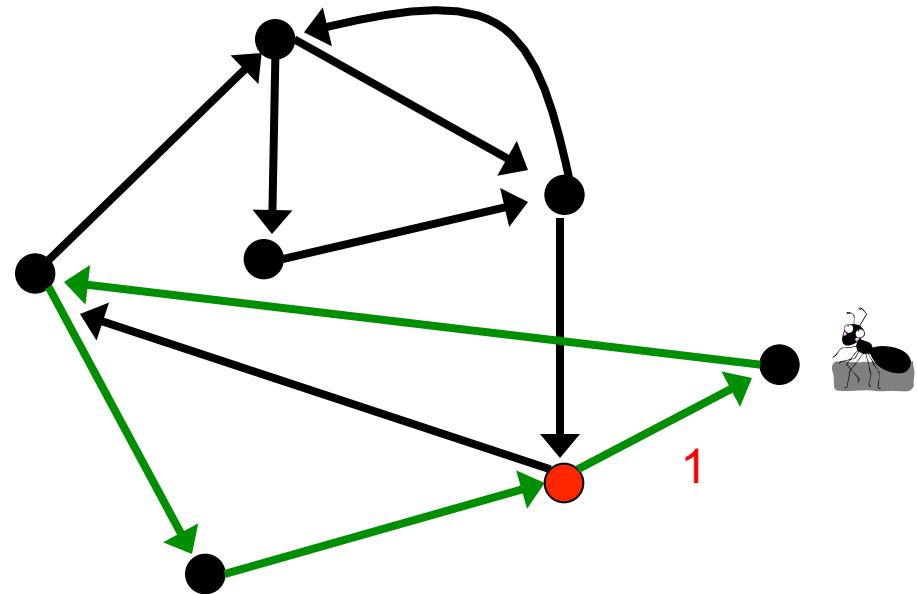
New Instructions for the Ant:

Starting at a **node** that has an unused edge, traverse the already constructed (green cycle) and return back to the starting node.



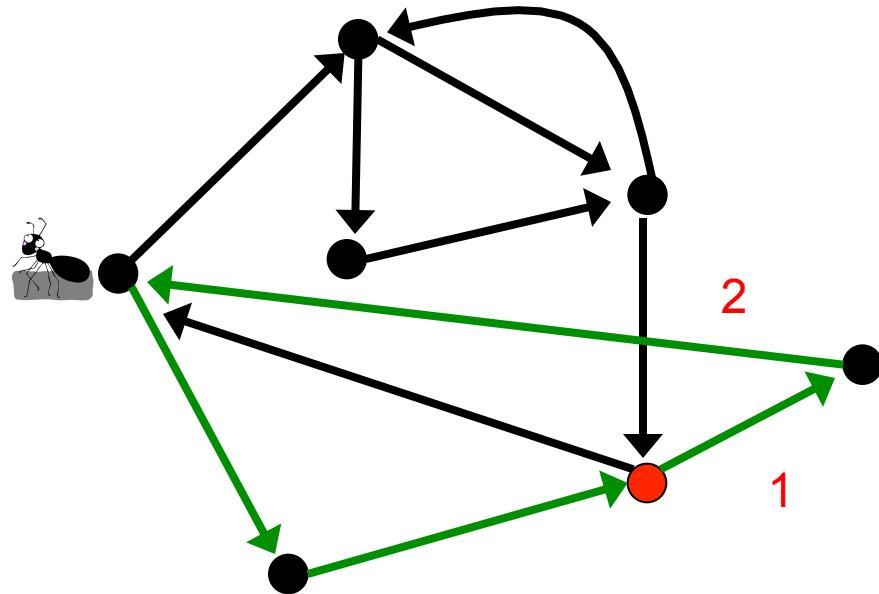
An Ant Traversing Previously Constructed Cycle

Starting at a node that has an unused edge, traverse the already constructed (green cycle) and return back to the starting node.



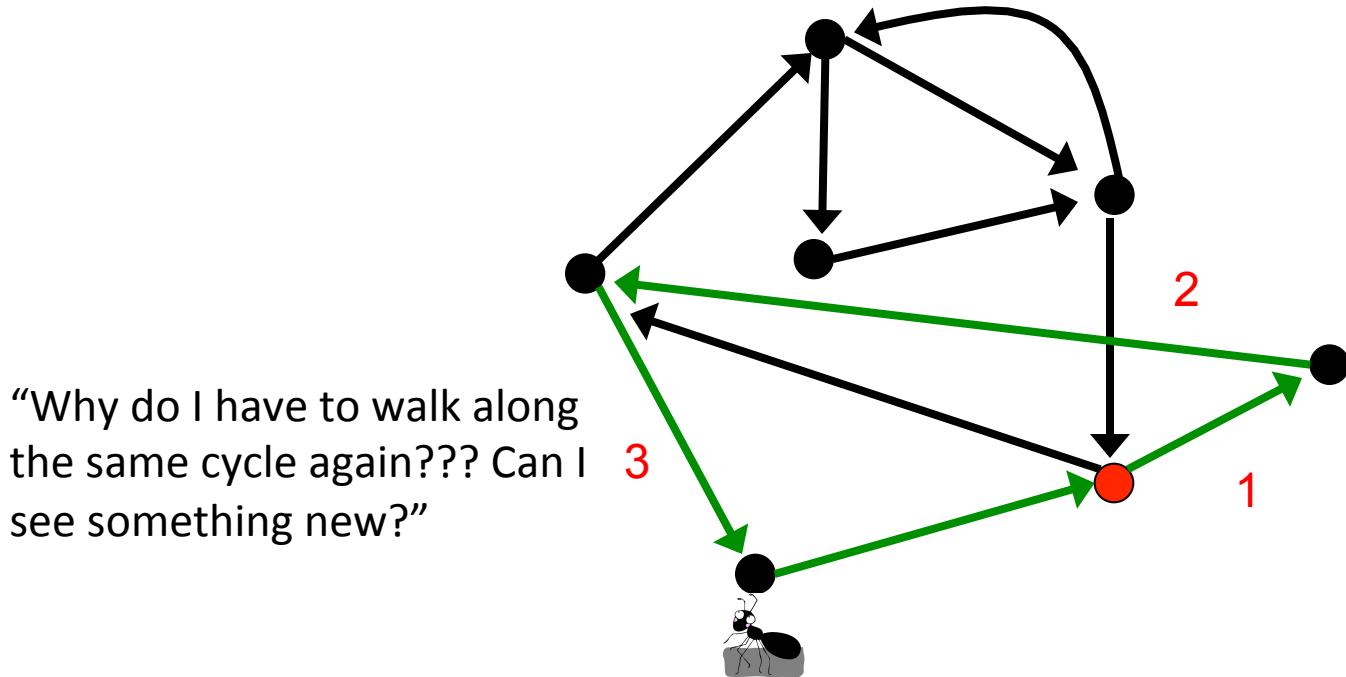
An Ant Traversing Previously Constructed Cycle

Starting at a node that has an unused edge, traverse the already constructed (green cycle) and return back to the starting node.



An Ant Traversing Previously Constructed Cycle

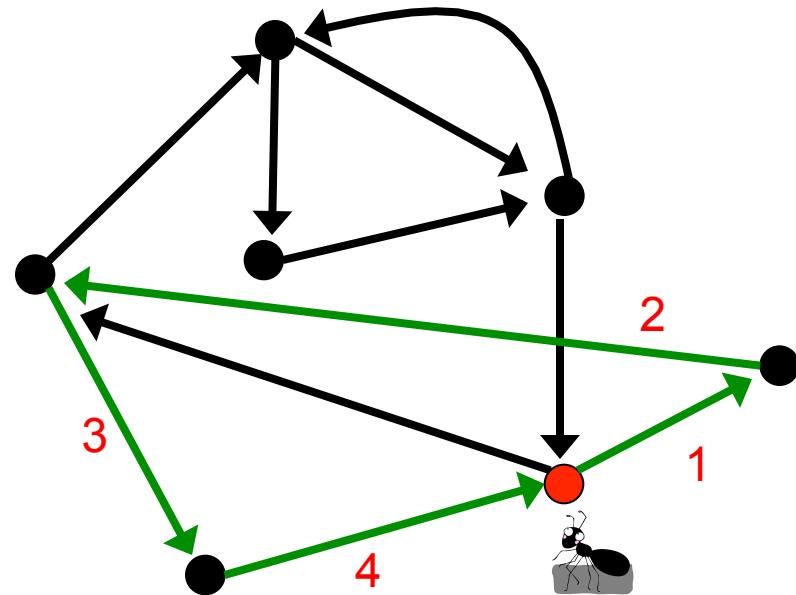
Starting at a node that has an unused edge, traverse the already constructed (green cycle) and return back to the starting node.



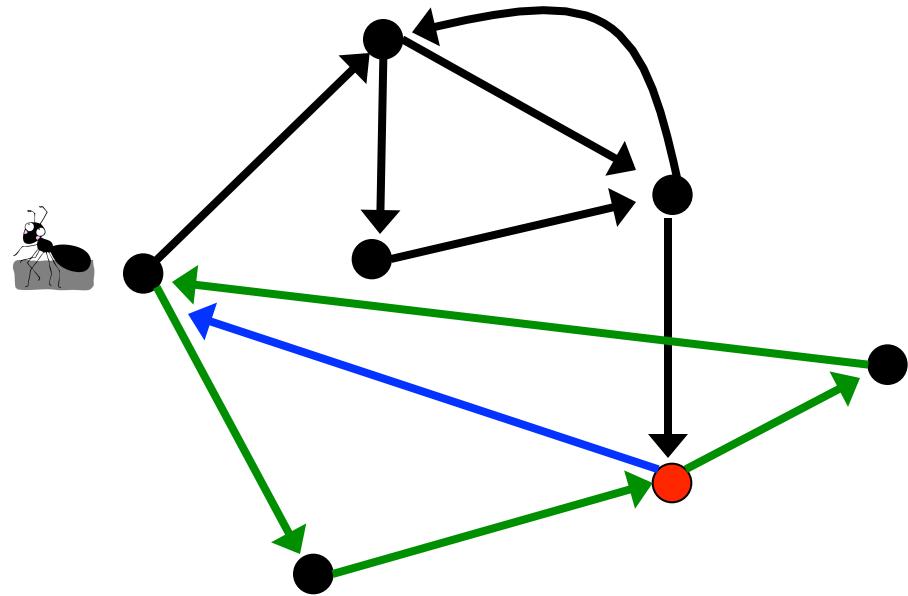
I Returned Back BUT... I Can Continue Walking!

Starting at a node that has an unused edge, traverse the already constructed (green cycle) and return back to the starting node.

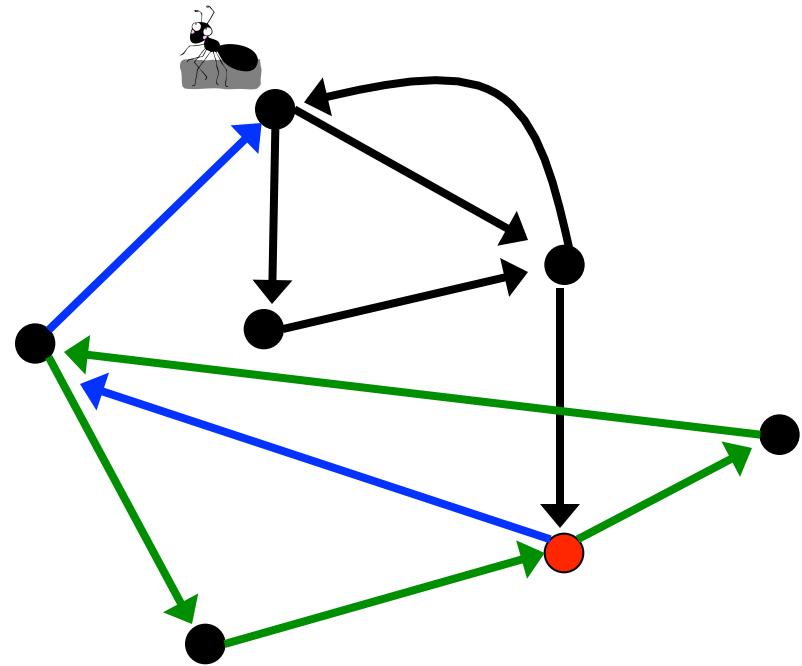
After completing the cycle, start random exploration of still untraversed edges in the graph.



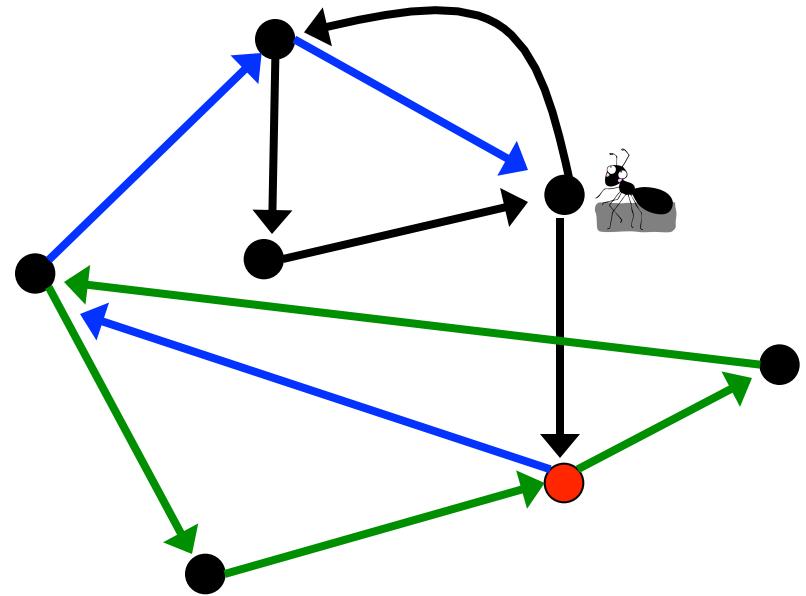
Enlarging the Previously Constructed Cycle



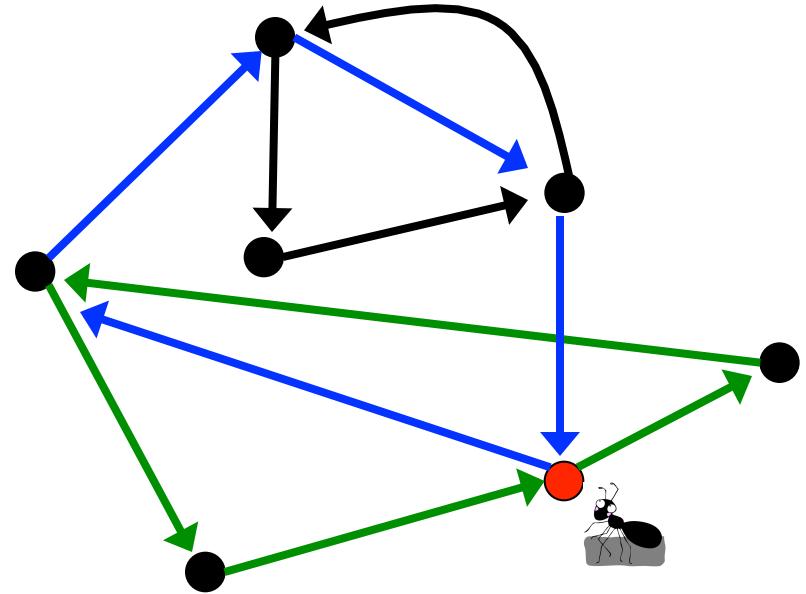
Enlarging the Previously Constructed Cycle



Enlarging the Previously Constructed Cycle



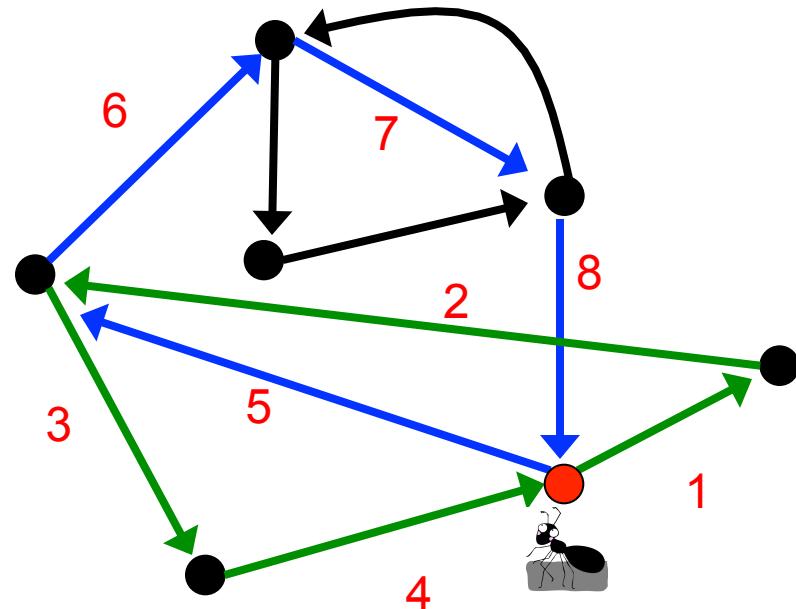
Enlarging the Previously Constructed Cycle



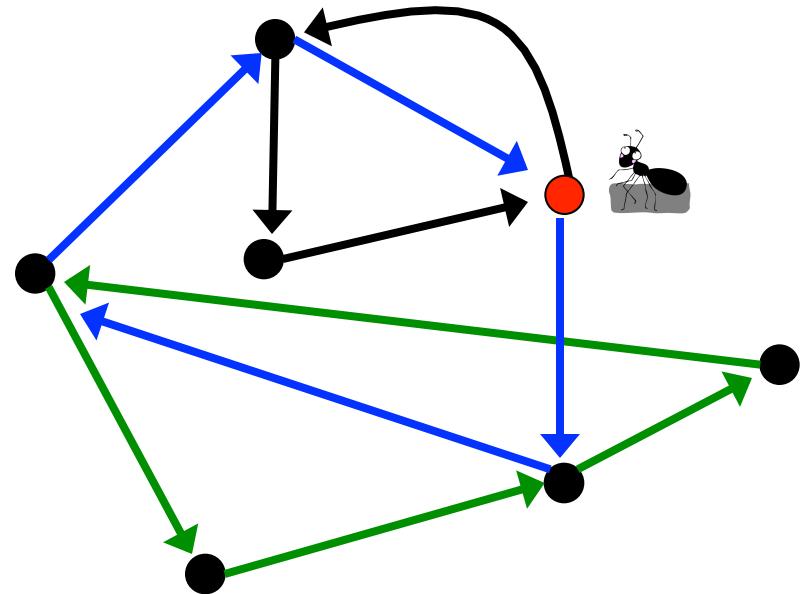
Stuck Again!

No Eulerian cycle yet... can we enlarge the green-blue cycle?

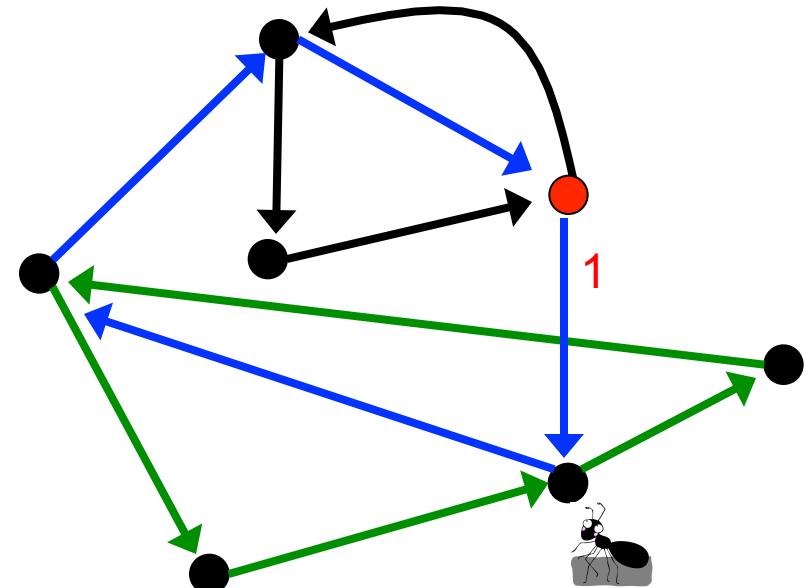
The ant should walk along the constructed cycle starting at yet another node. Which one?



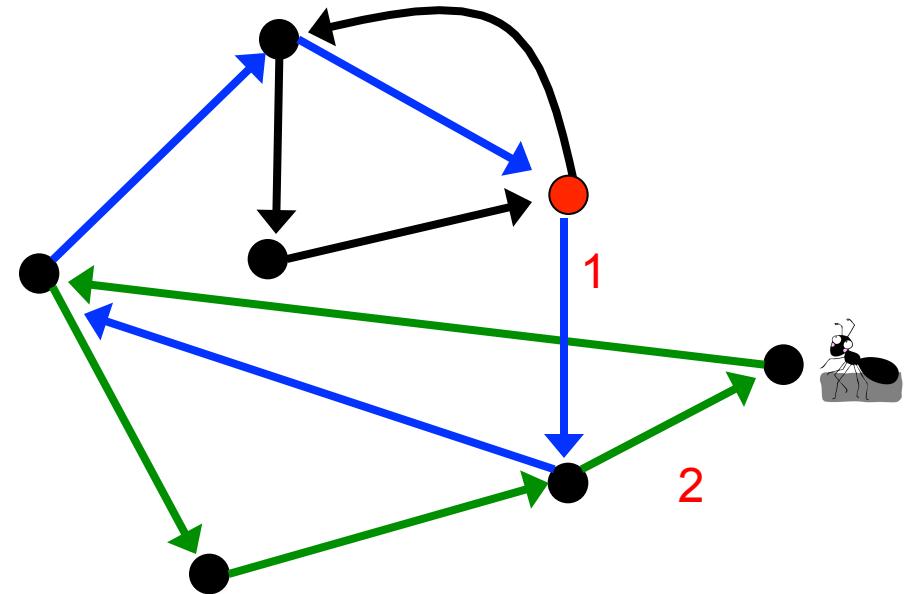
Starting at a New **Node**, Again...



Traversing the Previously Constructed Green-Blue Cycle

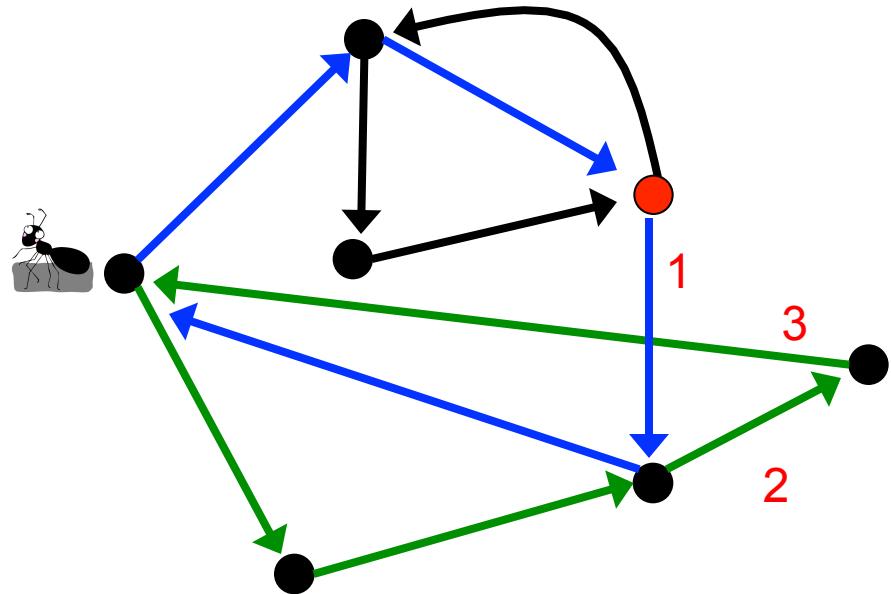


Traversing the Previously Constructed Green-Blue Cycle

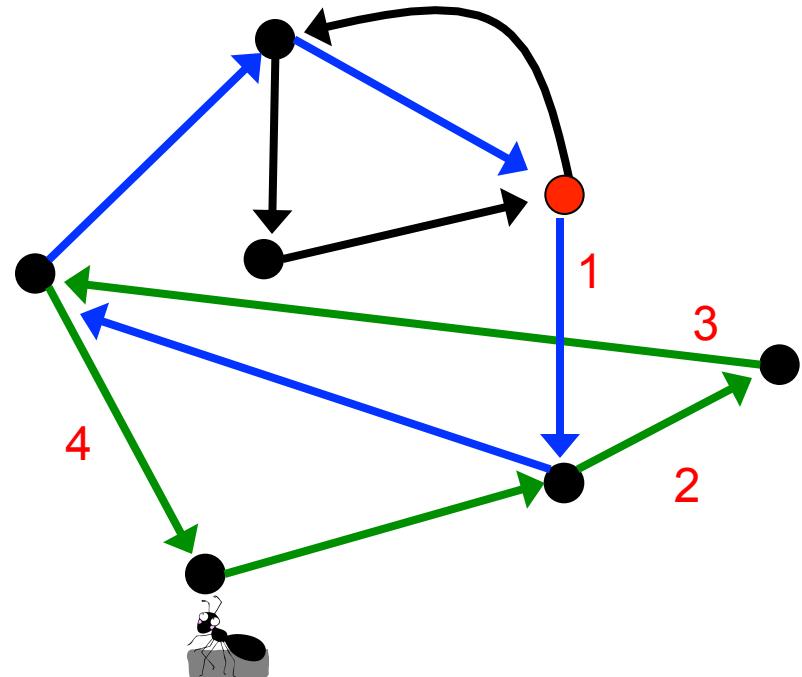


Traversing the Previously Constructed Green-Blue Cycle

"I hate to traverse the same cycle! What difference does it make where I start my walk???"

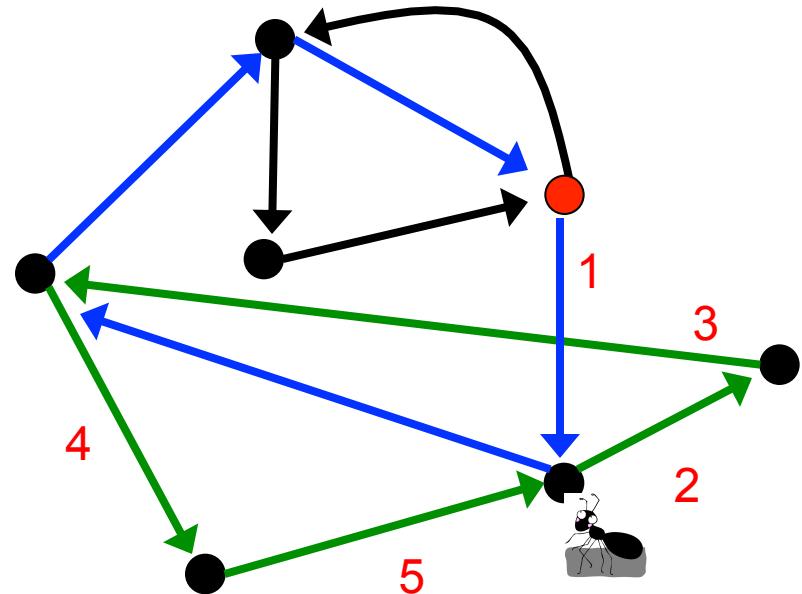


Traversing the Previously Constructed Green-Blue Cycle

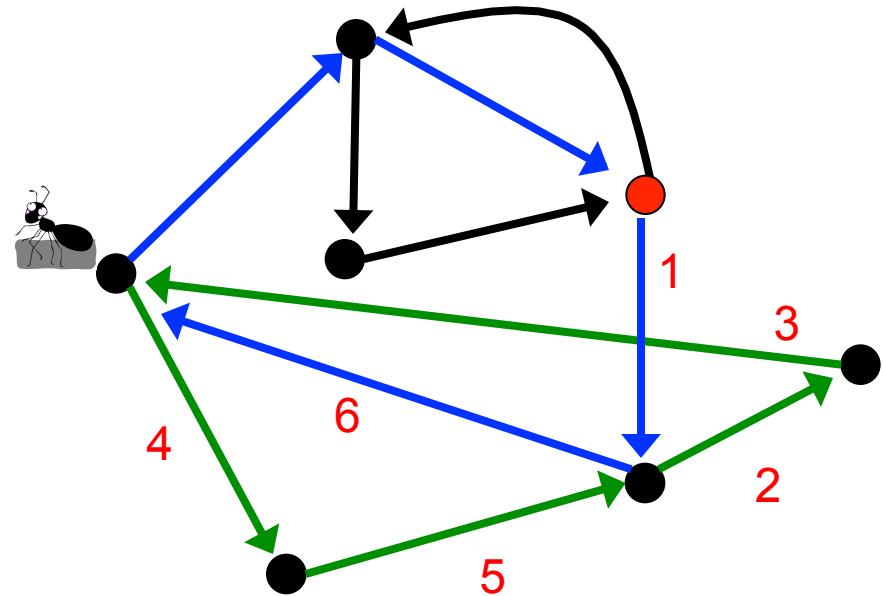


“These instructions are stupid...”

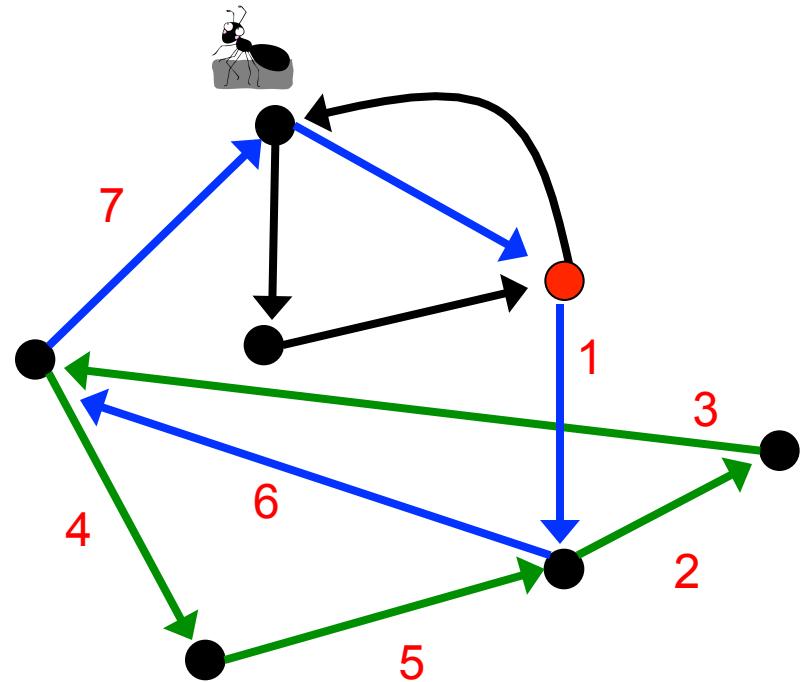
Traversing the Previously Constructed Green-Blue Cycle



Traversing the Previously Constructed Green-Blue Cycle

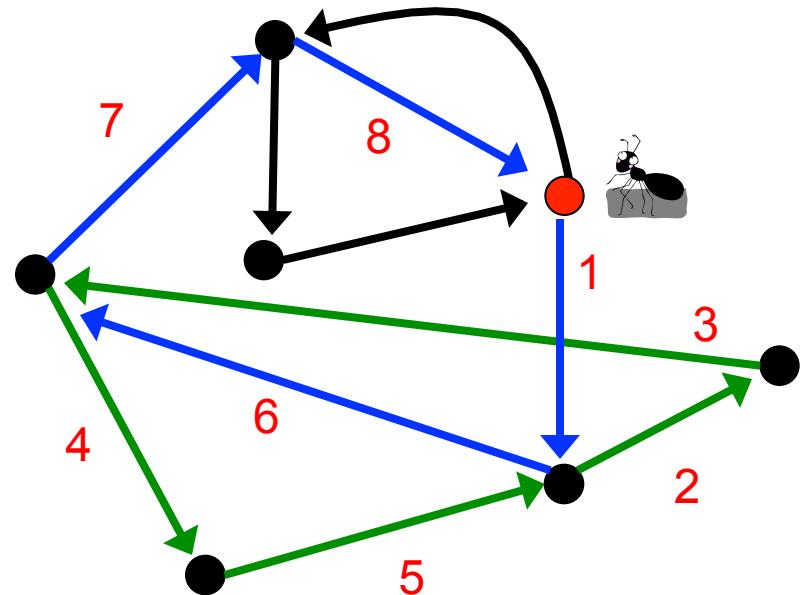


Traversing the Previously Constructed Green-Blue Cycle

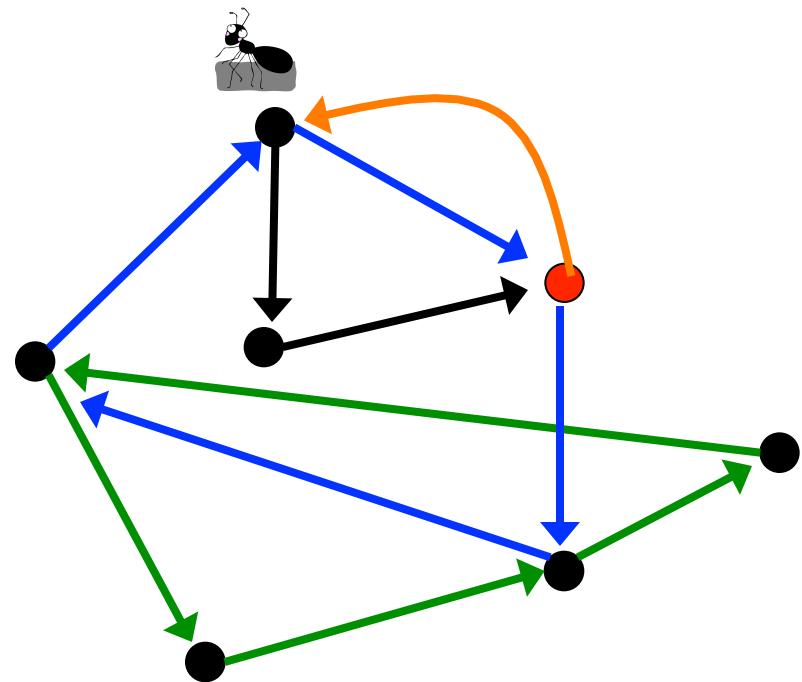


I Returned Back BUT... I Can Continue Walking!

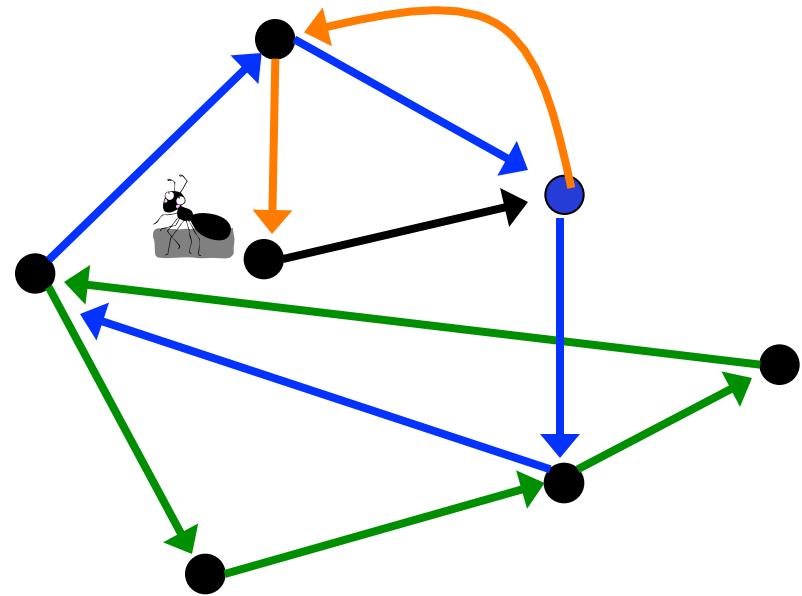
"Hmm, maybe these instructions were not that stupid..."



Enlarging the Green-Blue Cycle



Enlarging the Green-Blue Cycle



I Proved Euler's Theorem! Can I Go Home Please?

EulerianCycle(BalancedGraph)

form a *Cycle* by randomly walking in *BalancedGraph* (avoiding already visited edges)

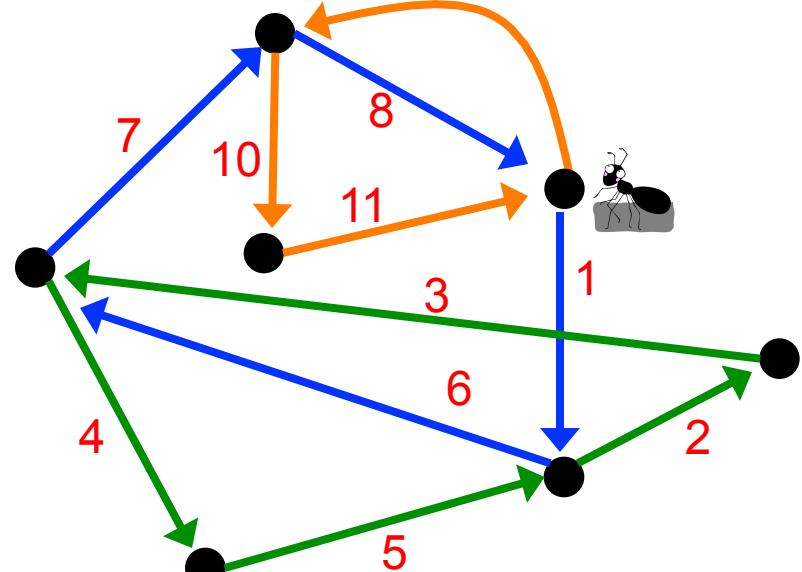
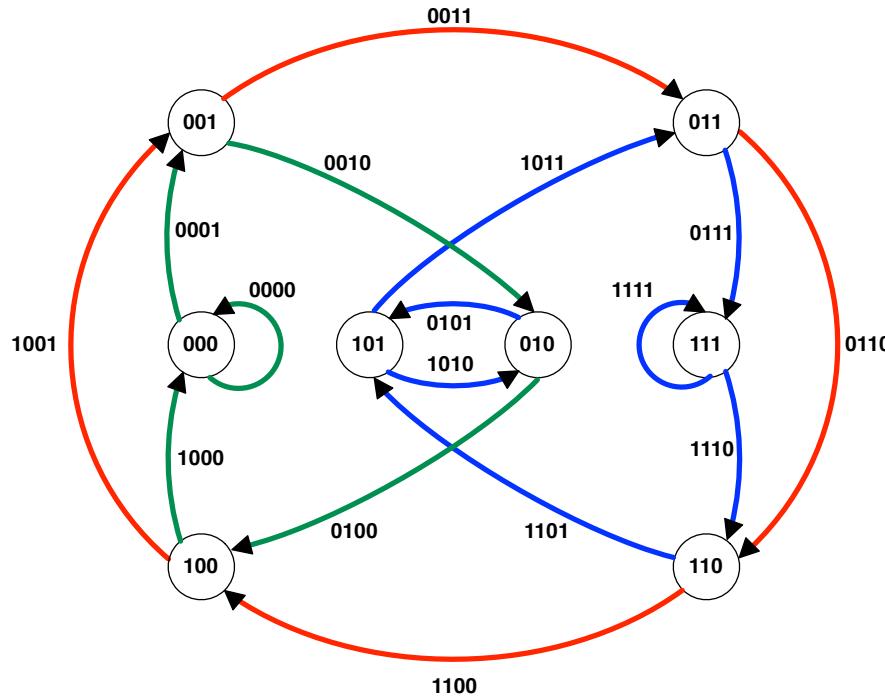
while *Cycle* is not Eulerian

 select a node *newStart* in *Cycle* with still unexplored outgoing edges

 form a *Cycle'* by traversing *Cycle* from *newStart* and randomly walking afterwards

Cycle \leftarrow *Cycle'*

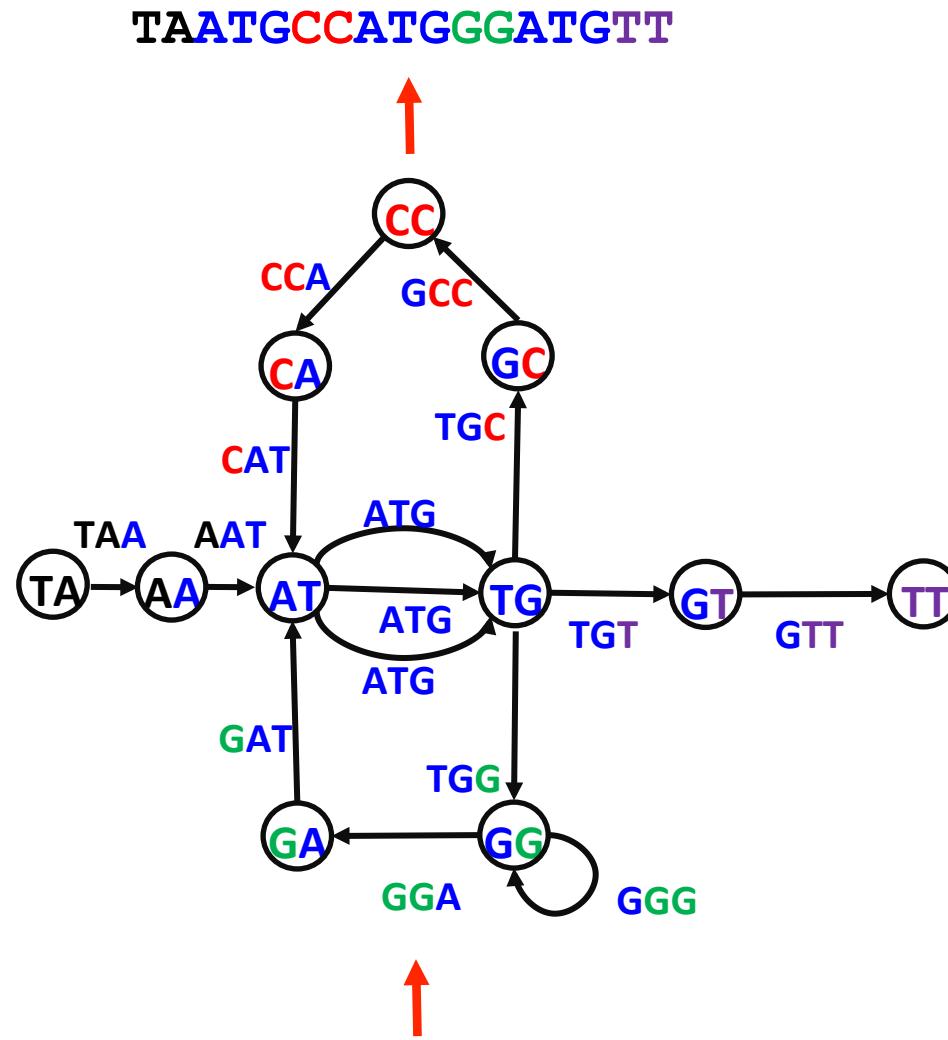
return *Cycle*



Outline

- What Is Genome Sequencing?
- Exploding Newspapers
- The String Reconstruction Problem
- String Reconstruction as a Hamiltonian Path Problem
- String Reconstruction as an Eulerian Path Problem
- Similar Problems with Different Fates
- De Bruijn Graphs
- Euler's Theorem
- **Assembling Read-Pairs**
- De Bruijn Graphs Face Harsh Realities of Assembly

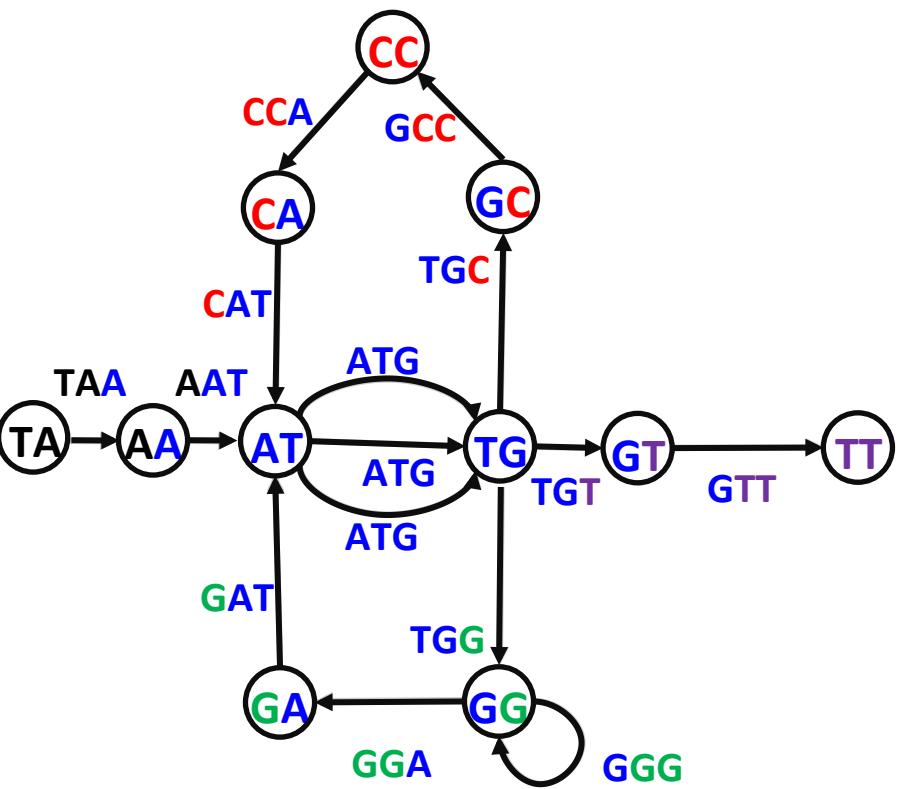
From Reads to de Bruijn Graph to Genome



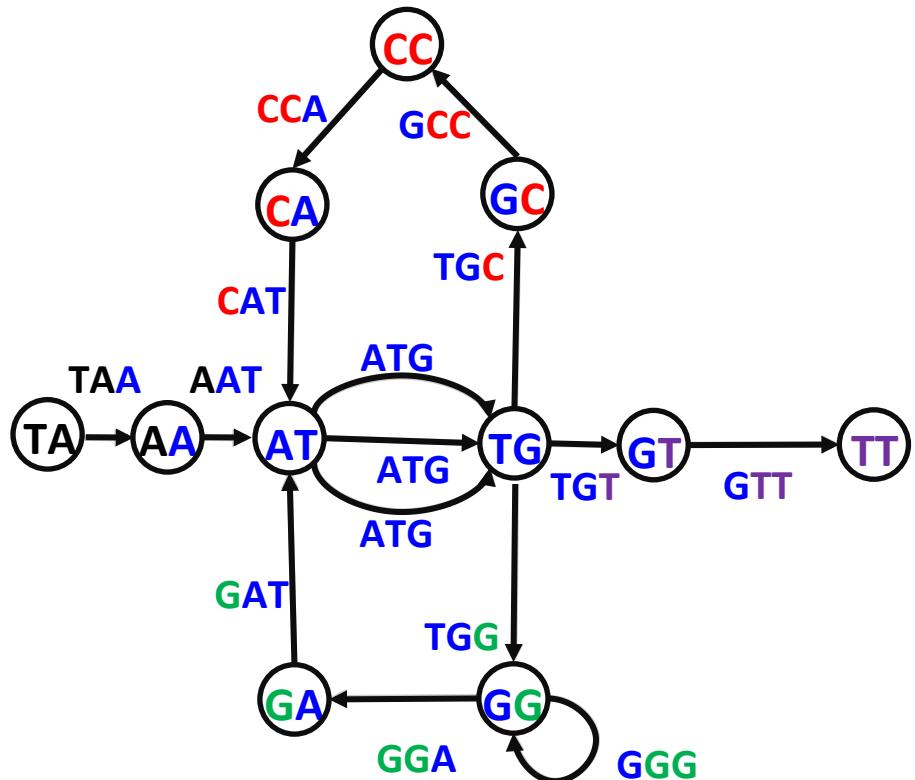
AAT ATG ATG ATG CAT CCA GAT GCC GGA GGG GTT TAA TGC TGG TGT

Multiple Eulerian Paths

TA~~ATGCCATGGGATGTT~~

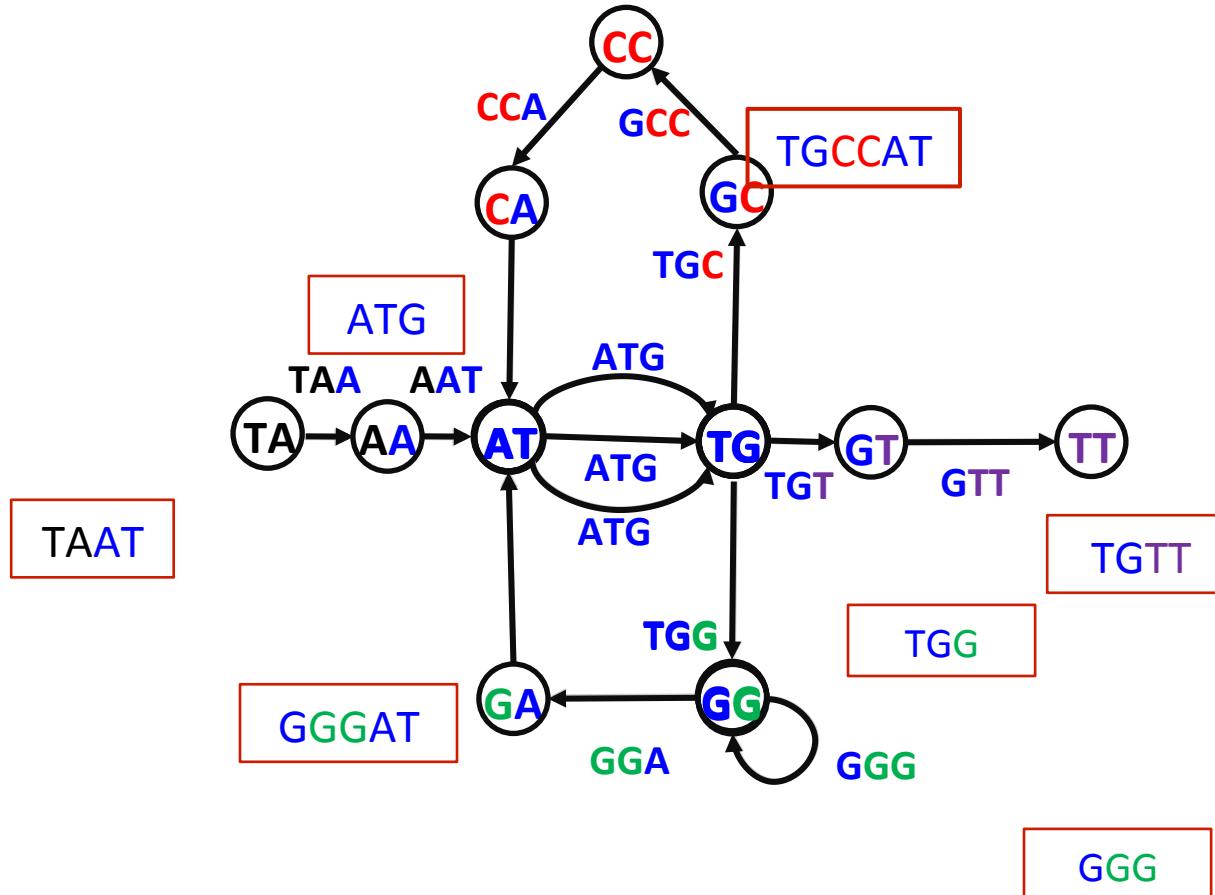


TA~~ATGGGATGCCATGTT~~



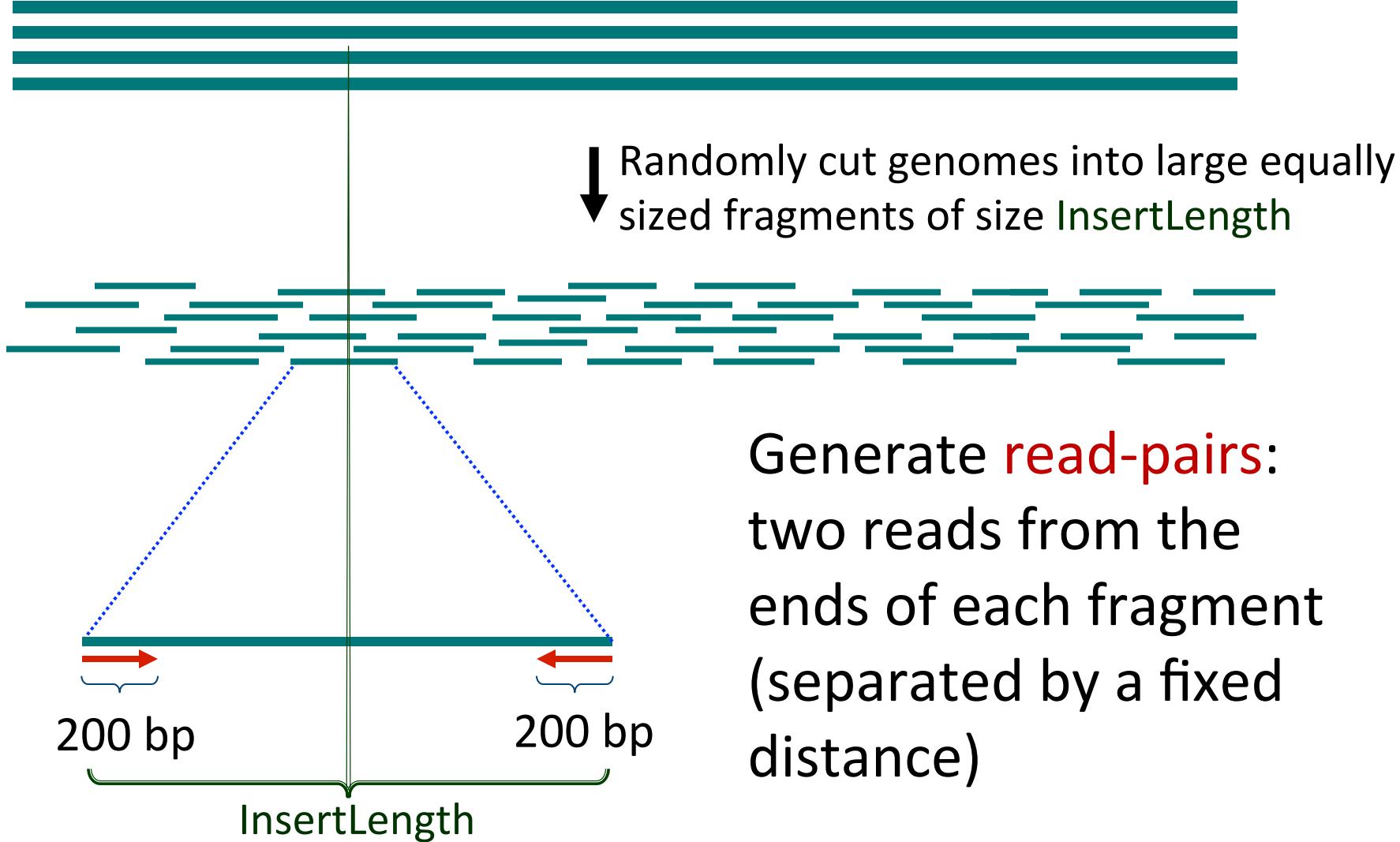
Breaking Genome into Contigs

TAATG**CC**CATGGGATGTT

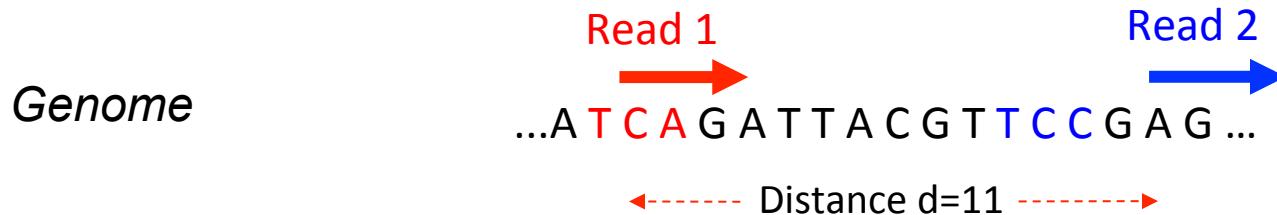


DNA Sequencing with Read-pairs

Multiple identical copies of genome



From k -mers to Paired k -mers



A paired k -mer is a pair of k -mers at a fixed distance d apart in Genome.
E.g. **TCA** and **TCC** are at distance $d=11$ apart.

Disclaimers:

1. In reality, **Read1** and **Read2** are typically sampled from different strands:
($\rightarrow \dots \dots \leftarrow$ rather than $\rightarrow \dots \dots \rightarrow$)
2. In reality, the distance d between reads is measured with errors.

What is PairedComposition(**T****A****ATGCC****CATGGG****ATGTT**)?

TAA **GCC**

paired 3-mer

What is PairedComposition(**TAA****ATGCC****CATGGG****ATGTT**)?

TAA **GCC**
AAT **CCA**
ATG **CAT**
TGC **ATG**
GCC **TGG**
CCA **GGG**
CAT **GGA**
ATG **GAT**
TGG **ATG**
GGG **TGT**
GGA **GTT**

Representing a paired 3-mer **TAA** **GCC** as a 2-line expression:

TAA
GCC

TAA **AAT** **ATG** **TGC** **GCC** **CCA** **CAT** **ATG** **TGG** **GGG** **GGA**
GCC **CCA** **CAT** **ATG** **TGG** **GGG** **GGA** **GAT** **ATG** **TGT** **GTT**

PairedComposition(TAATGCCATGGGATGTT)

TAA GCC
AAT CCA
ATG CAT
TGC ATG
GCC TGG
CCA GGG
CAT GGA
ATG GAT
TGG ATG
GGG TGT
GGA GTT

TAA GCC	AAT CCA	ATG CAT	TGC ATG	GCC TGG	CCA GGG	CAT GGA	ATG GAT	TGG ATG	GGG TGT	GGA GTT
AAT CCA	ATG CAT	ATG GAT	CAT GGA	CCA GGG	GCC TGG	GGA GTT	GGG TGT	TAA GCC	TGC ATG	TGG ATG

Representing PairedComposition in lexicographic order

String Reconstruction from Read-Pairs Problem

String Reconstruction from Read-Pairs Problem. Reconstruct a string from its paired k -mers.

- **Input.** A collection of paired k -mers.
- **Output.** A string $Text$ such that $\text{PairedComposition}(Text)$ is equal to the collection of paired k -mers.

How Would de Bruijn Assemble Paired k -mers?

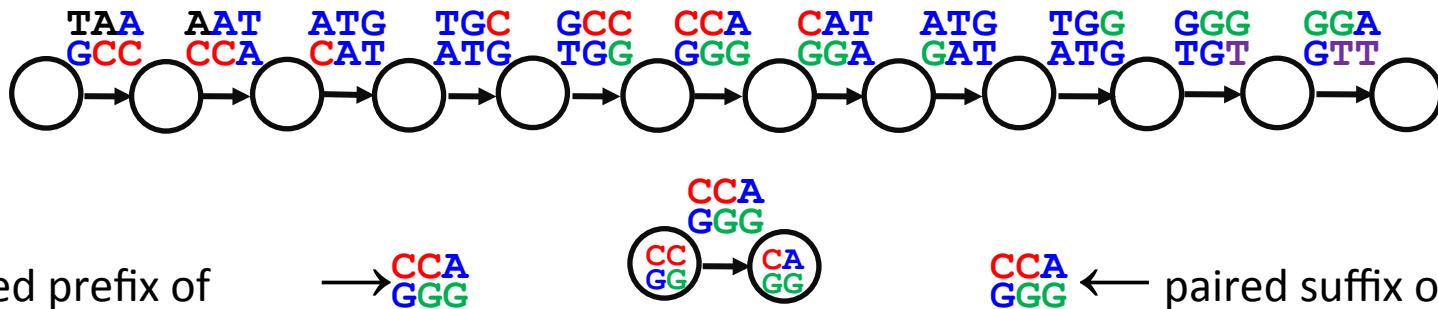


Paired de Bruijn Graphs

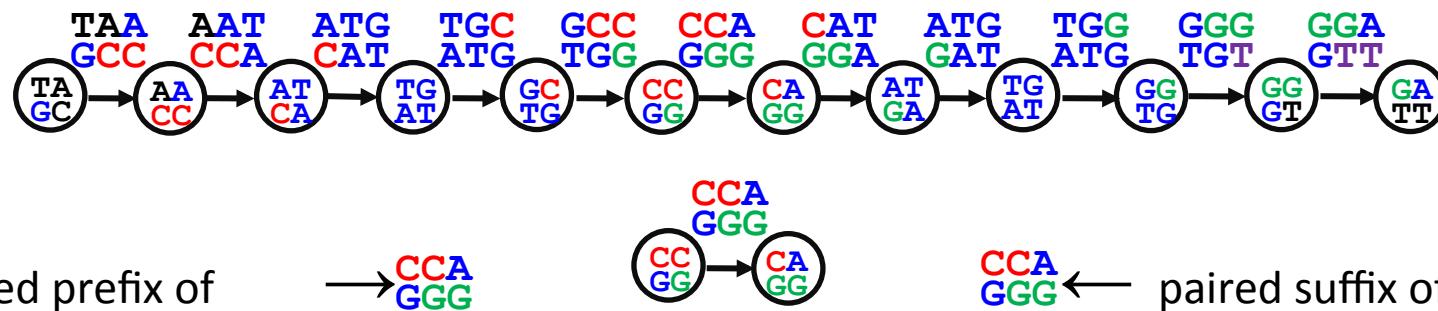


Representing Genome **TAATGCCATGGGATGTT** as a Path

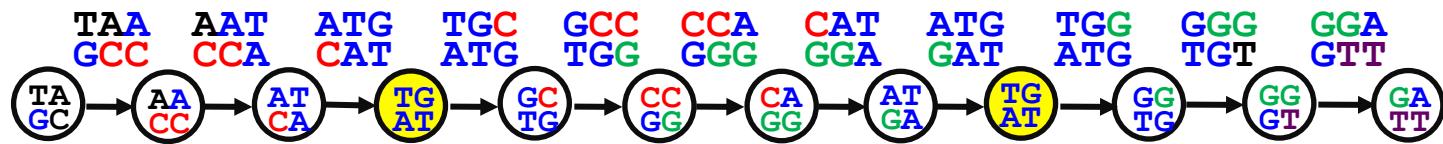
TAA GCC
AAT CCA
ATG CAT
TGC ATG
GCC TGG
CCA GGG
CAT GGA
ATG GAT
TGG ATG
GGG TGT
GGA GTT



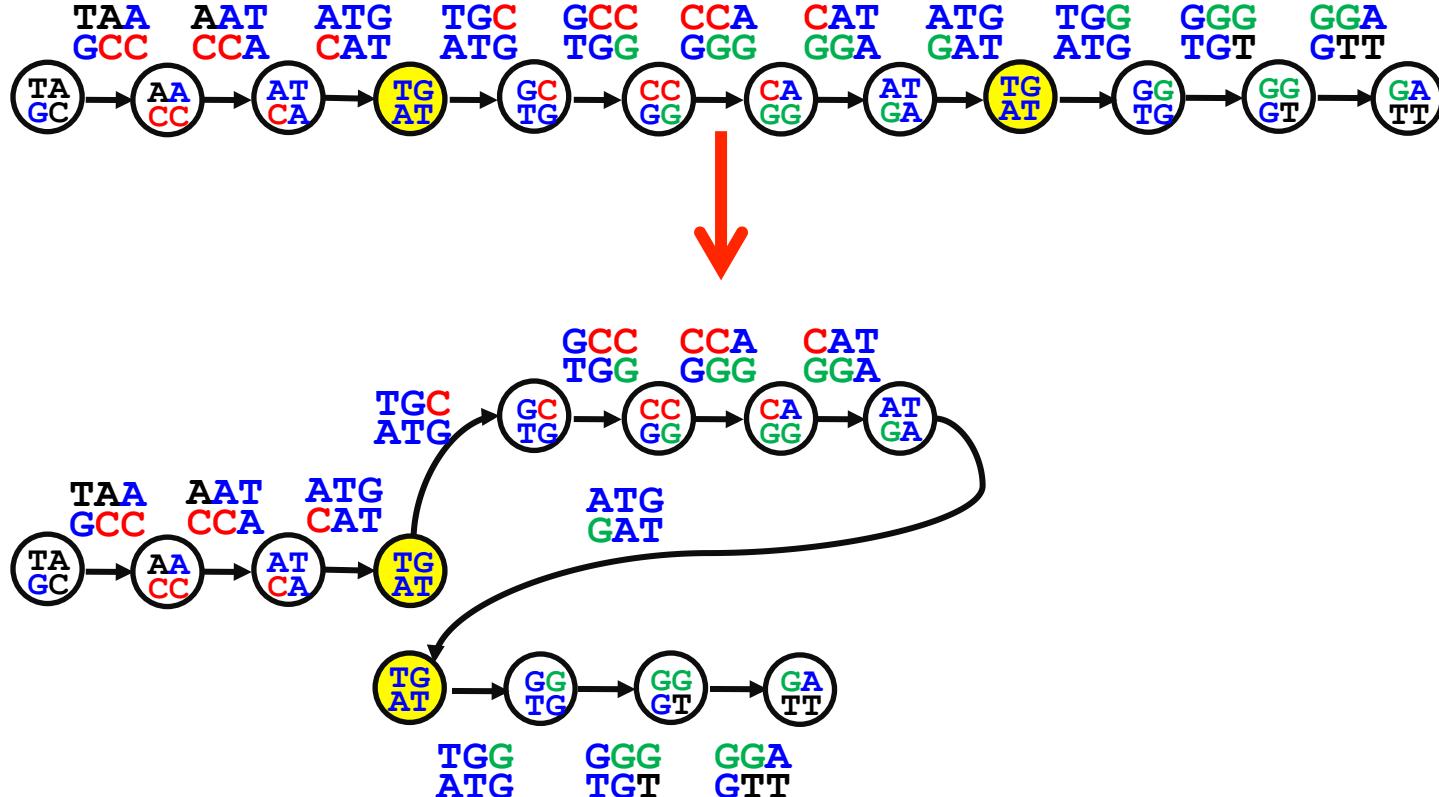
Labeling Nodes by Paired Prefixes and Suffixes



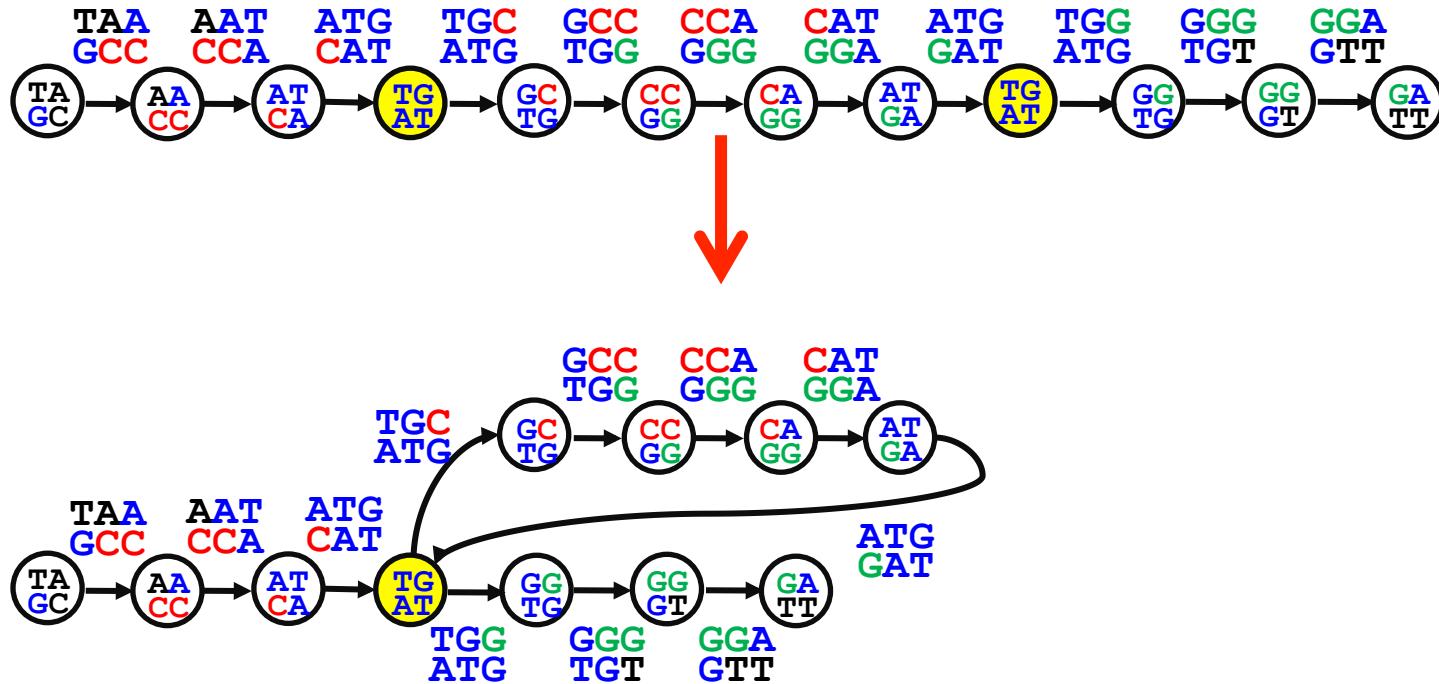
Glue nodes with identical labels



Glue nodes with identical labels



Glue nodes with identical labels



Paired de Bruijn Graph from the Genome

Constructing Paired de Bruijn Graph from paired k-mers

TAA
GCC

ATG
CAT

GCC
TGG

CAT
GGA

TGG
ATG

GGA
GTT

AAT
CCA

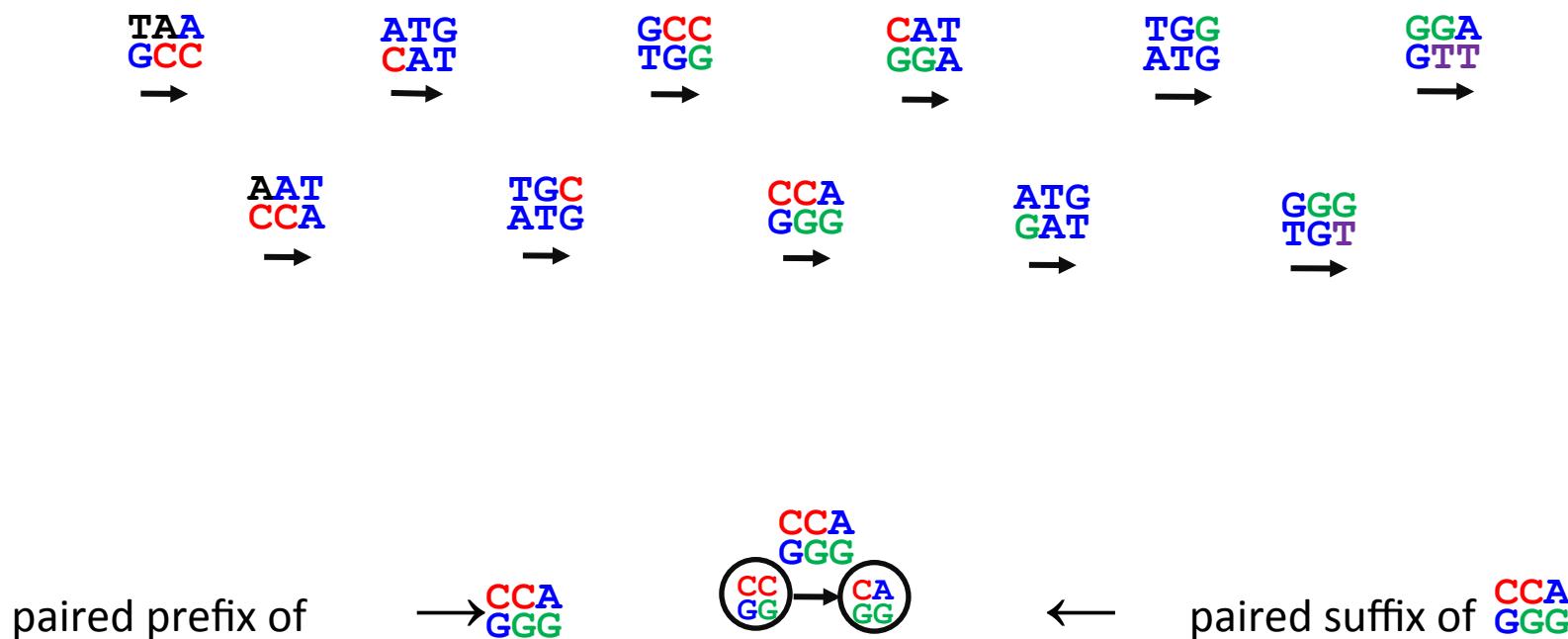
TGC
ATG

CCA
GGG

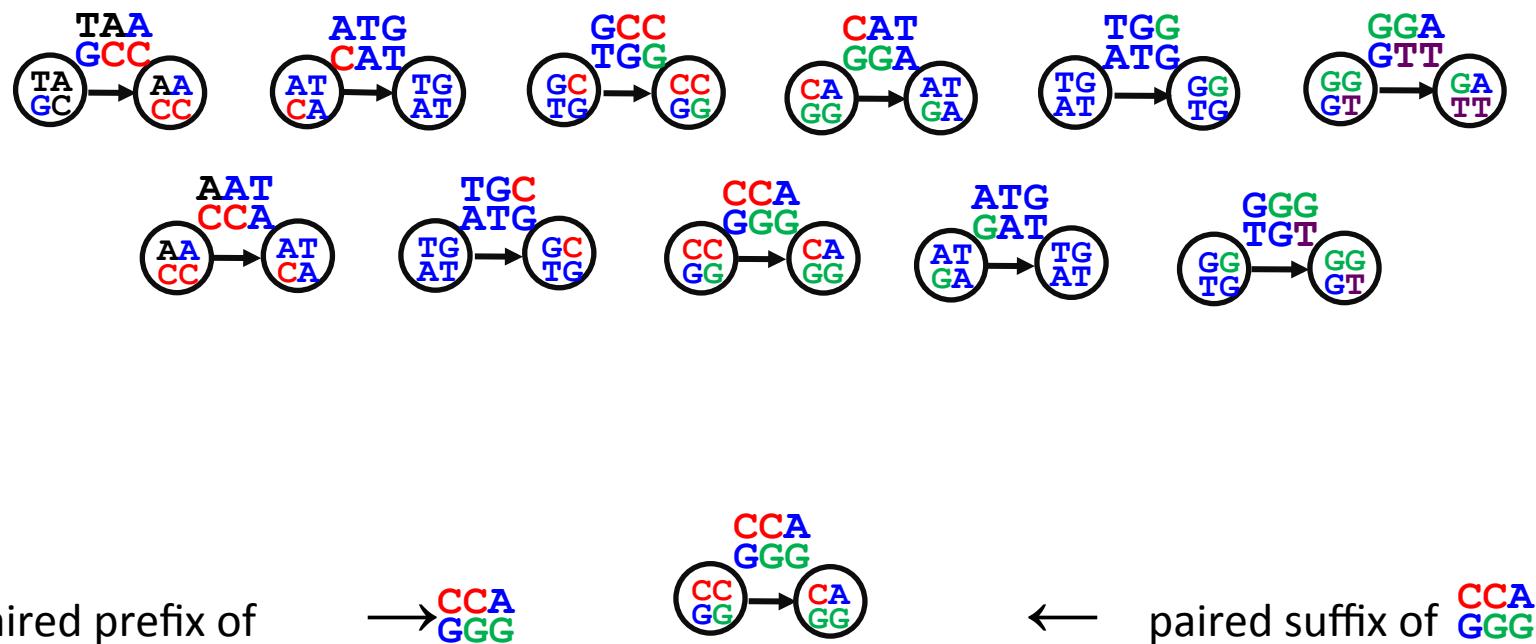
ATG
GAT

GGG
TGT

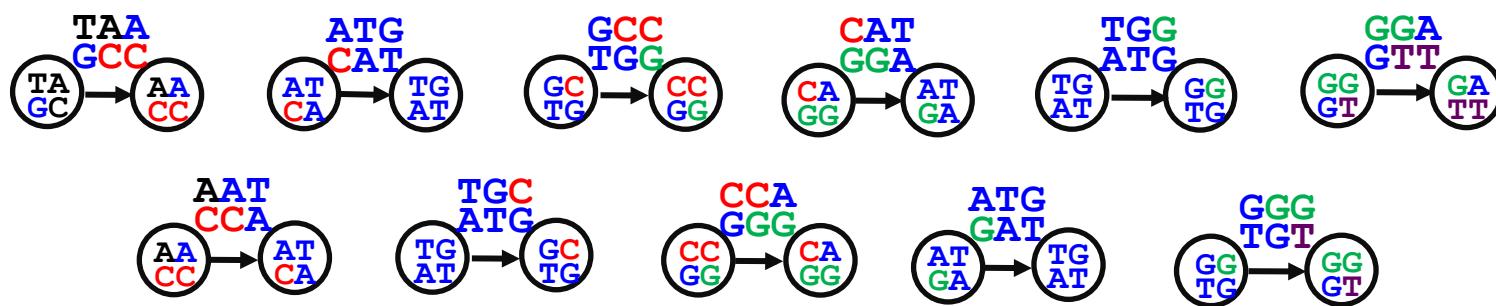
Constructing Paired de Bruijn Graph from paired k-mers



Constructing Paired de Bruijn Graph

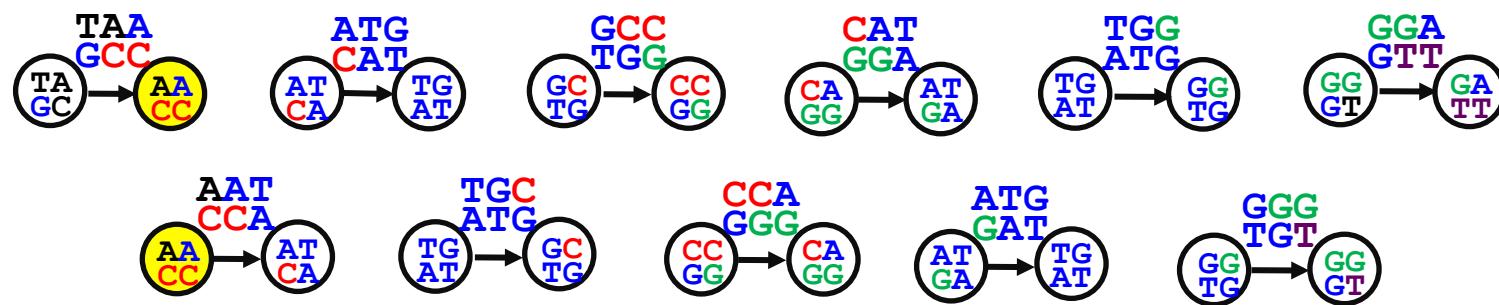


Constructing Paired de Bruijn Graph

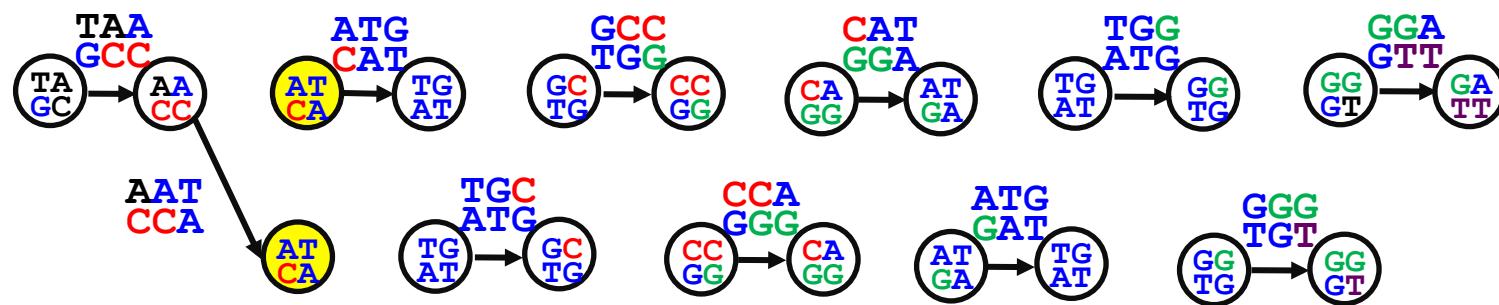


- **Paired de Bruijn graph for a collection of paired k -mers:**
 - Represent every paired k -mer as an edge between its paired prefix and paired suffix.
 - Glue **ALL** nodes with identical labels.

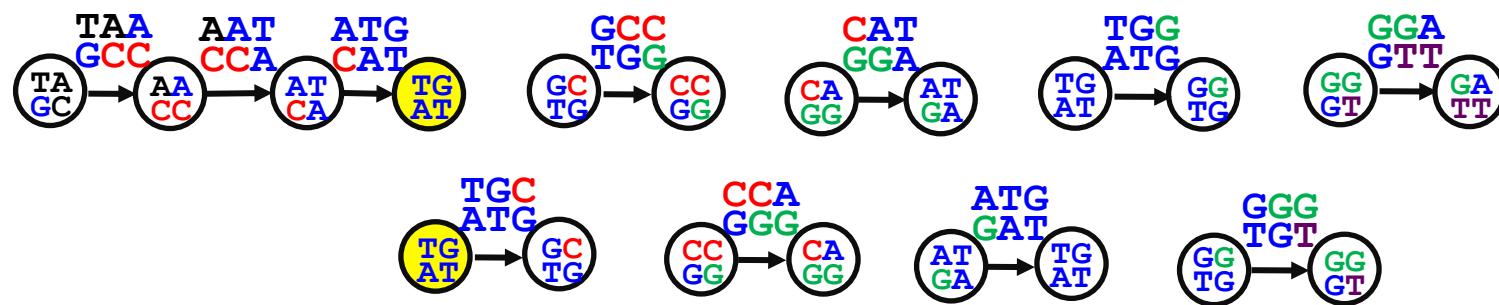
Constructing Paired de Bruijn Graph



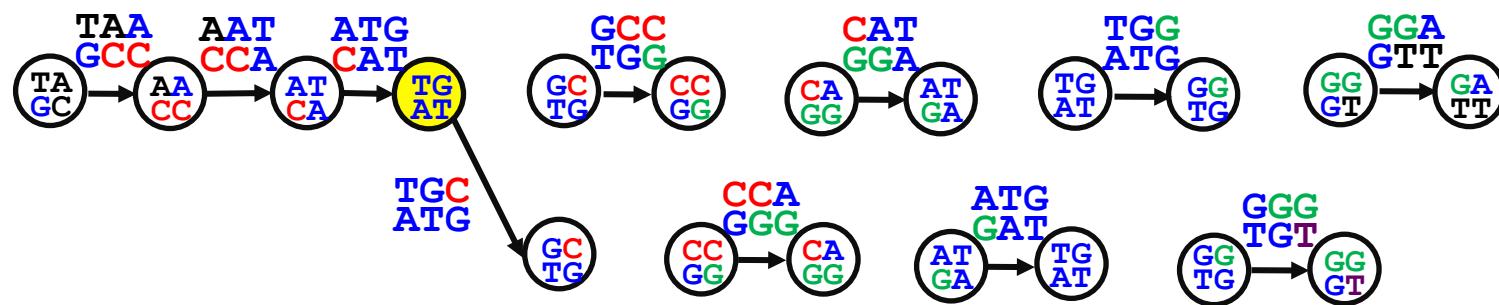
Constructing Paired de Bruijn Graph



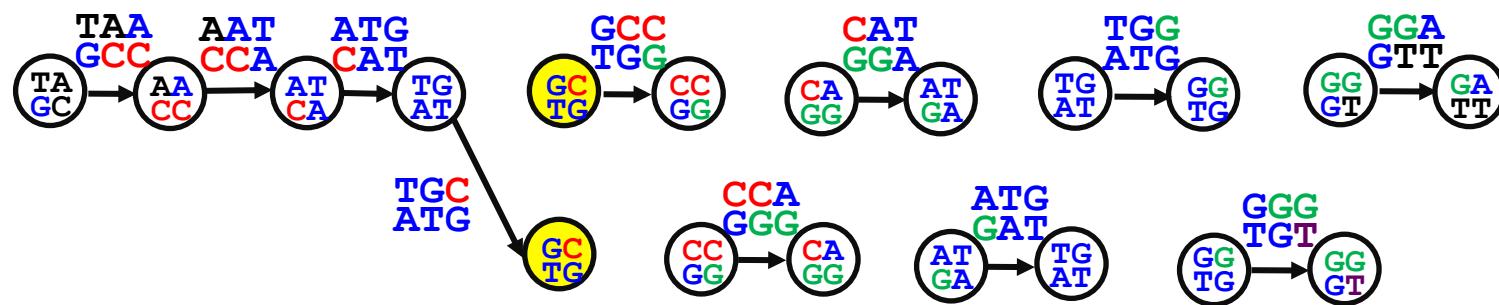
Constructing Paired de Bruijn Graph



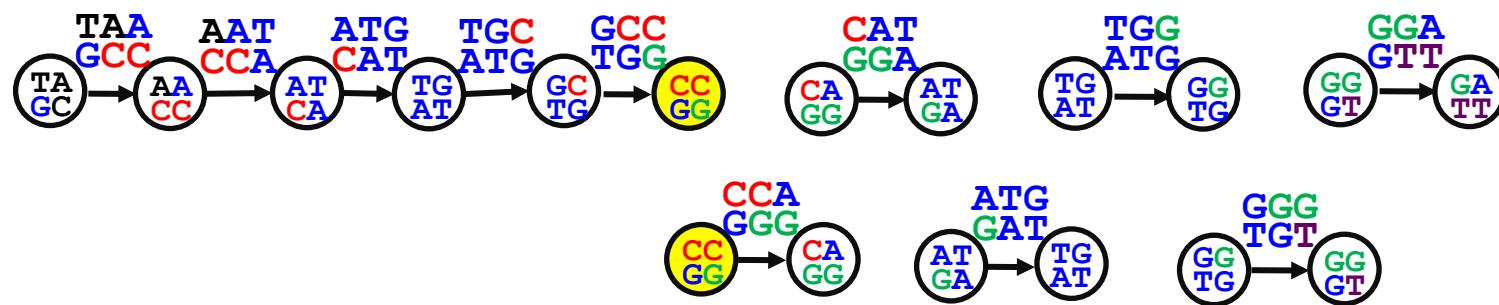
Constructing Paired de Bruijn Graph



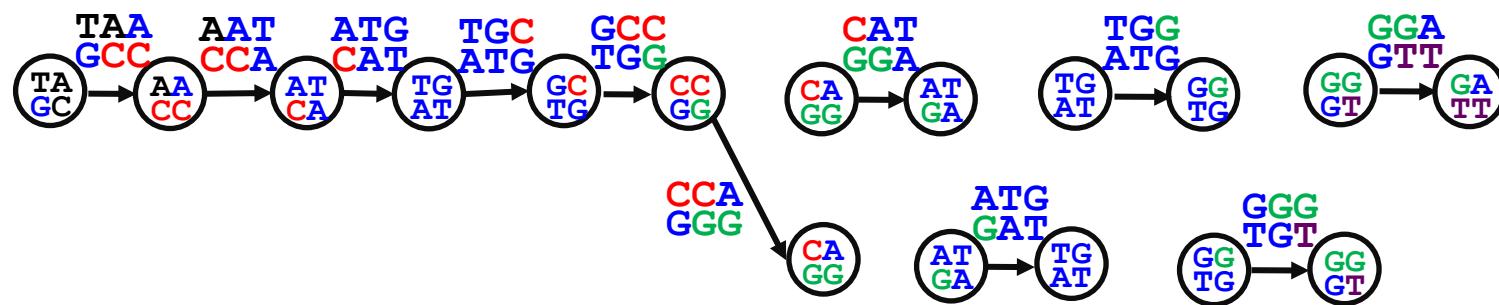
Constructing Paired de Bruijn Graph



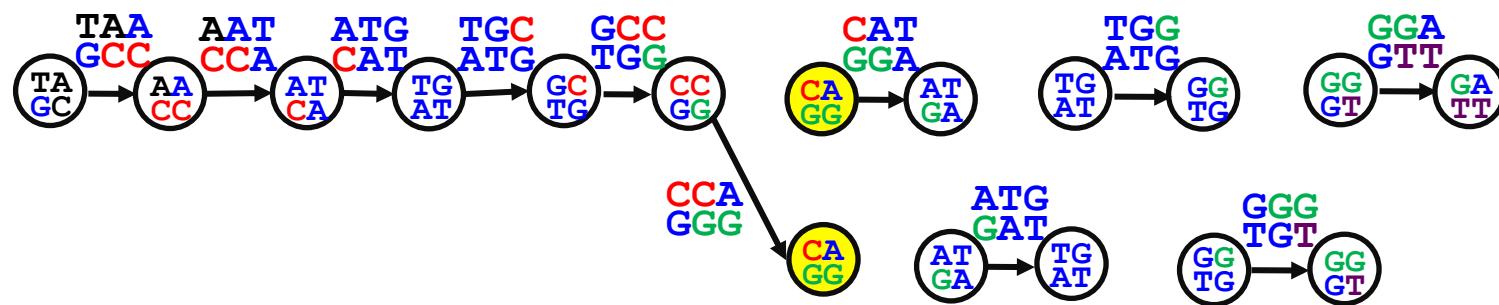
Constructing Paired de Bruijn Graph



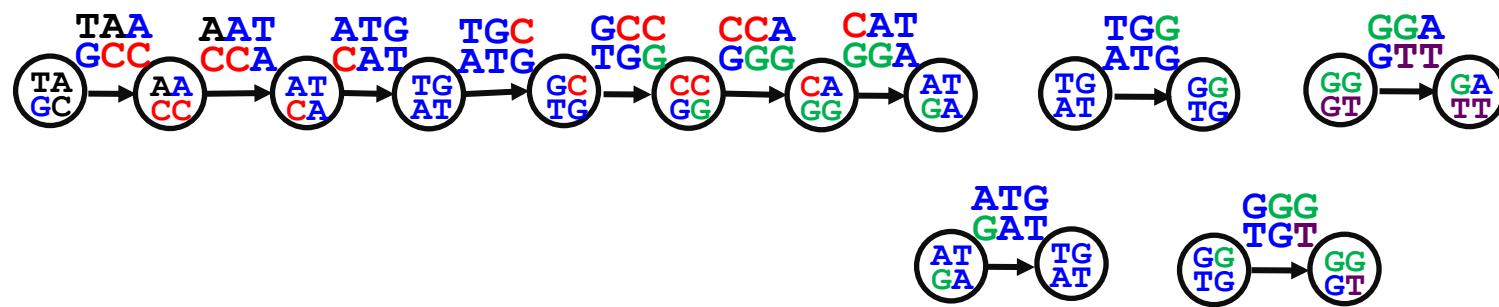
Constructing Paired de Bruijn Graph



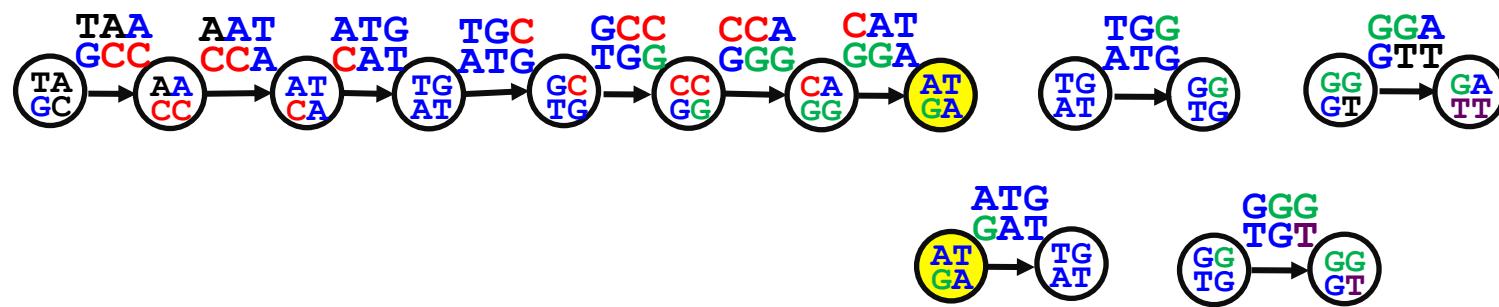
Constructing Paired de Bruijn Graph



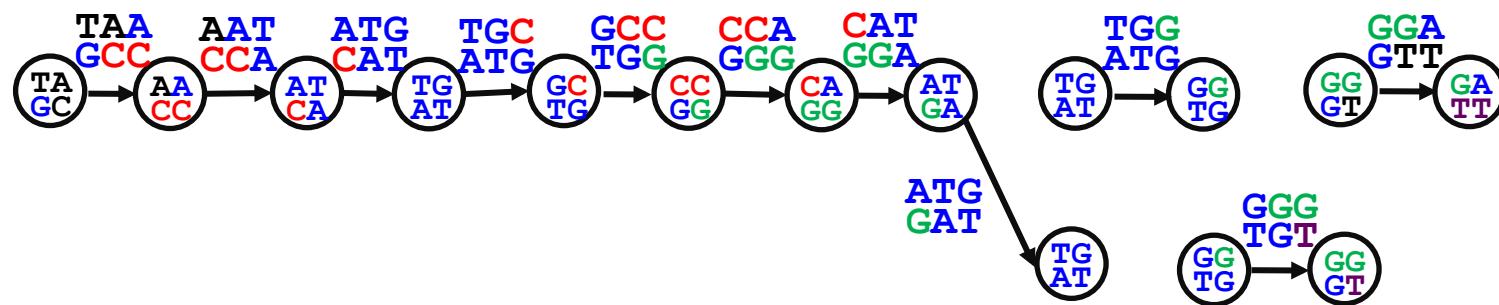
Constructing Paired de Bruijn Graph



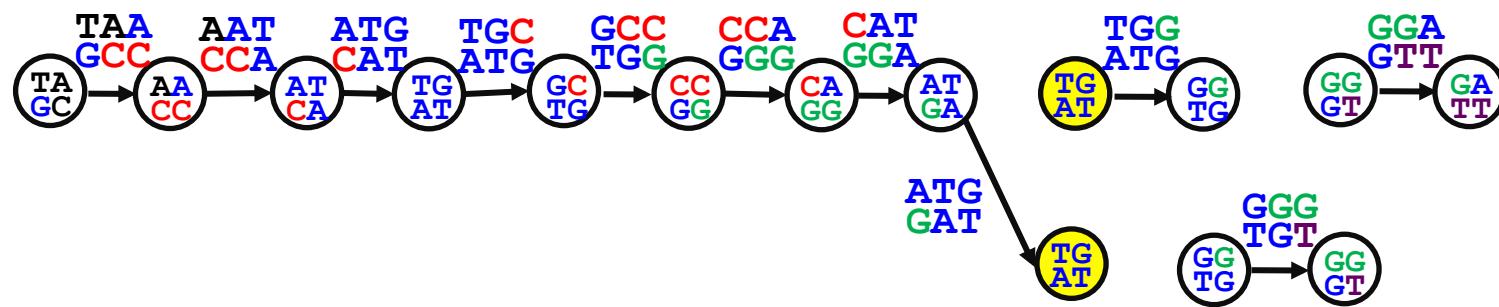
Constructing Paired de Bruijn Graph



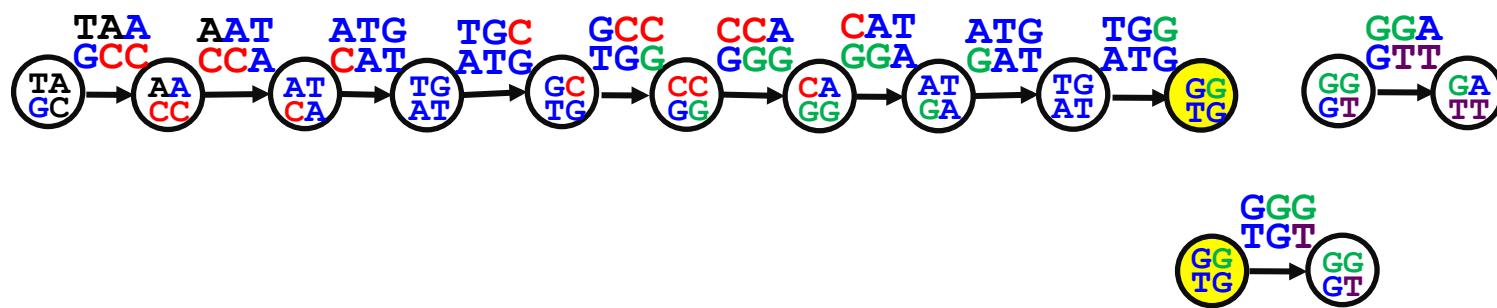
Constructing Paired de Bruijn Graph



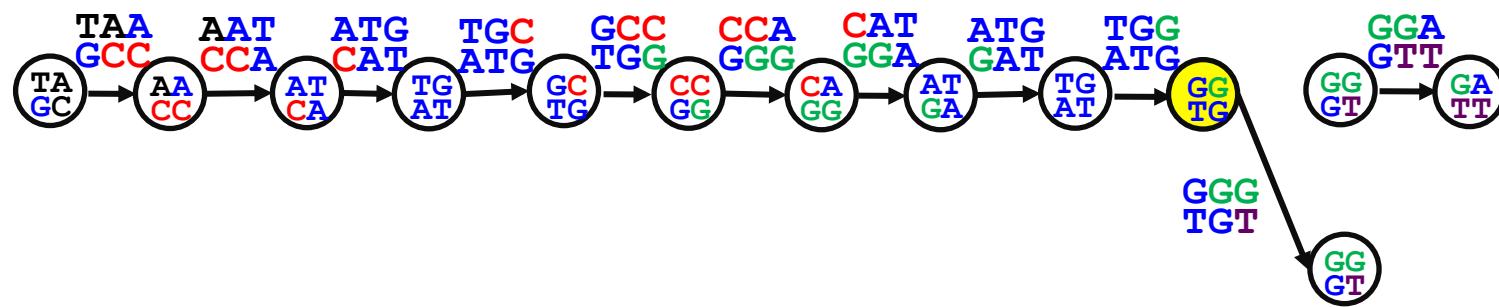
Constructing Paired de Bruijn Graph



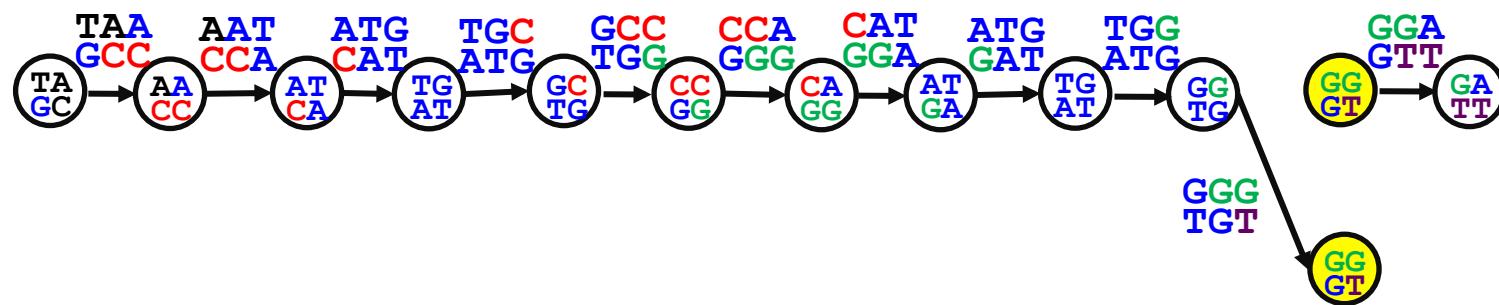
Constructing Paired de Bruijn Graph



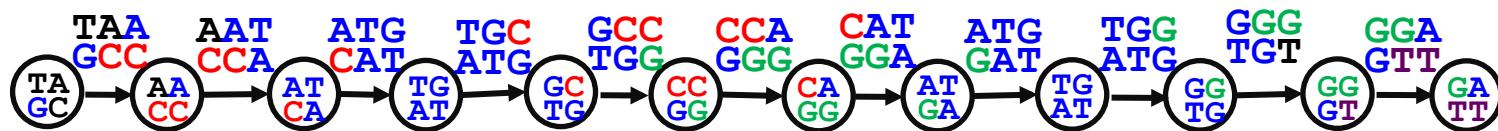
Constructing Paired de Bruijn Graph



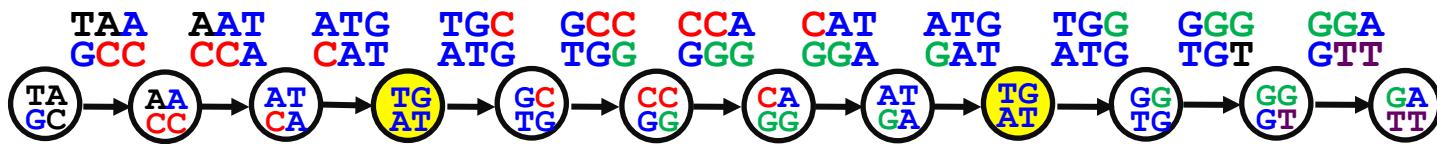
Constructing Paired de Bruijn Graph



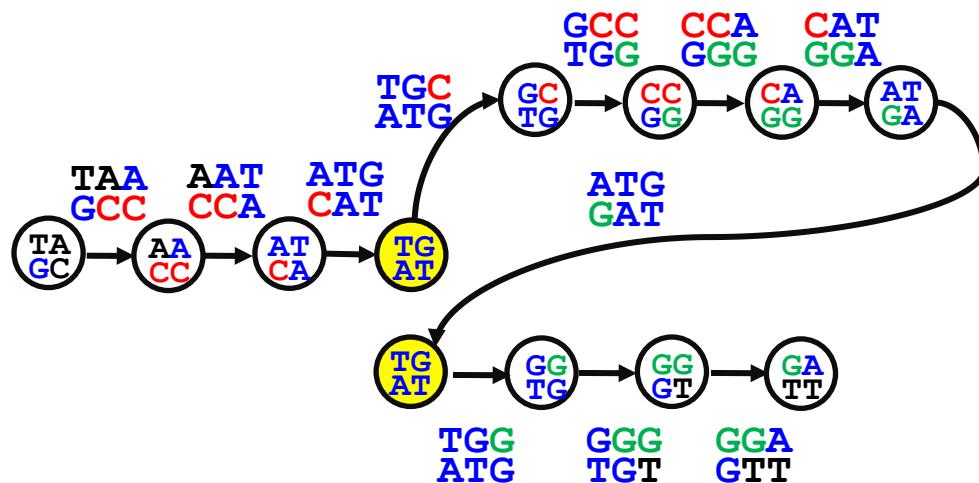
Constructing Paired de Bruijn Graph



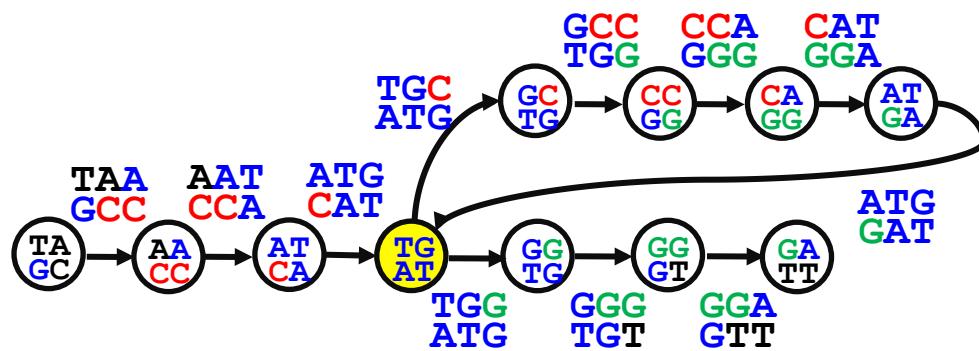
We Are Not Done with Gluing Yet



Constructing Paired de Bruijn Graph



Constructing Paired de Bruijn Graph



Paired de Bruijn Graph from read-pairs

Paired de Bruijn Graphs



- **Paired de Bruijn graph for a collection of paired k -mers:**
 - Represent every paired k -mer as an edge between its paired prefix and paired suffix.
 - Glue **ALL** nodes with identical labels.

Which Graph Represents a Better Assembly?

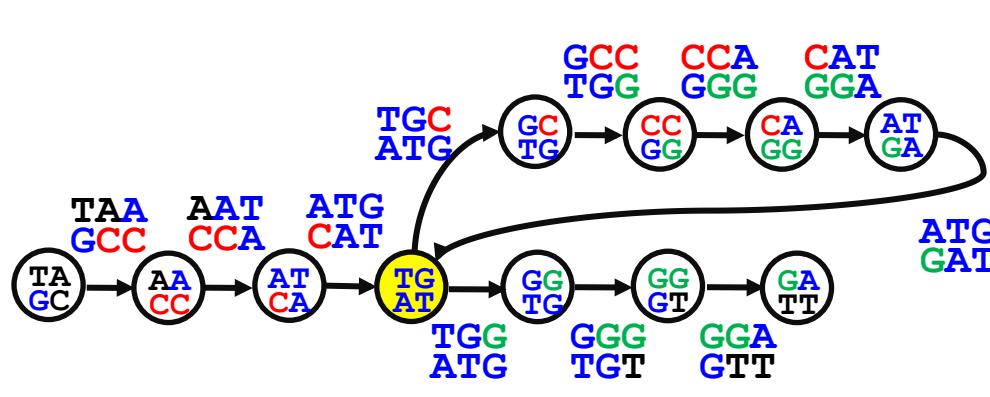
Unique genome reconstruction

Multiple genome reconstructions

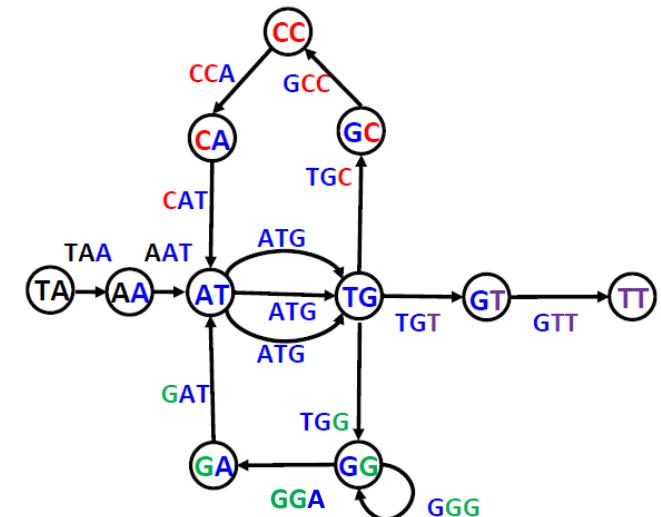
TAATGCCATGGGATGTT

TAATGCCATGGGATGTT

TAATGGGATGCCATGTT



Paired de Bruijn Graph



De Bruijn Graph