

# Genome Assembly

Lecture 22  
Oct 24, 2016

# Announcements

**Final Project:** The final project will consist of an oral presentation and written report (e.g., the methods section) related to an assembly project. Projects must incorporate an implementation of the computational techniques we've learned about. The final project will be worth 100 points (75 written/25 oral). Oral presentations will occur during the last 3 days of class. Written reports will be due on the last day of class. More details will be provided later in the semester.

# Announcements

<http://oyster-river-protocol.readthedocs.io/en/v2/>

<https://www.ebi.ac.uk/ena/>

Probably Transcriptome  
Illumina, PacBio, Nanopore  
Error Correct  
Trim  
Assembly  
Filter  
Annotate

# Announcements

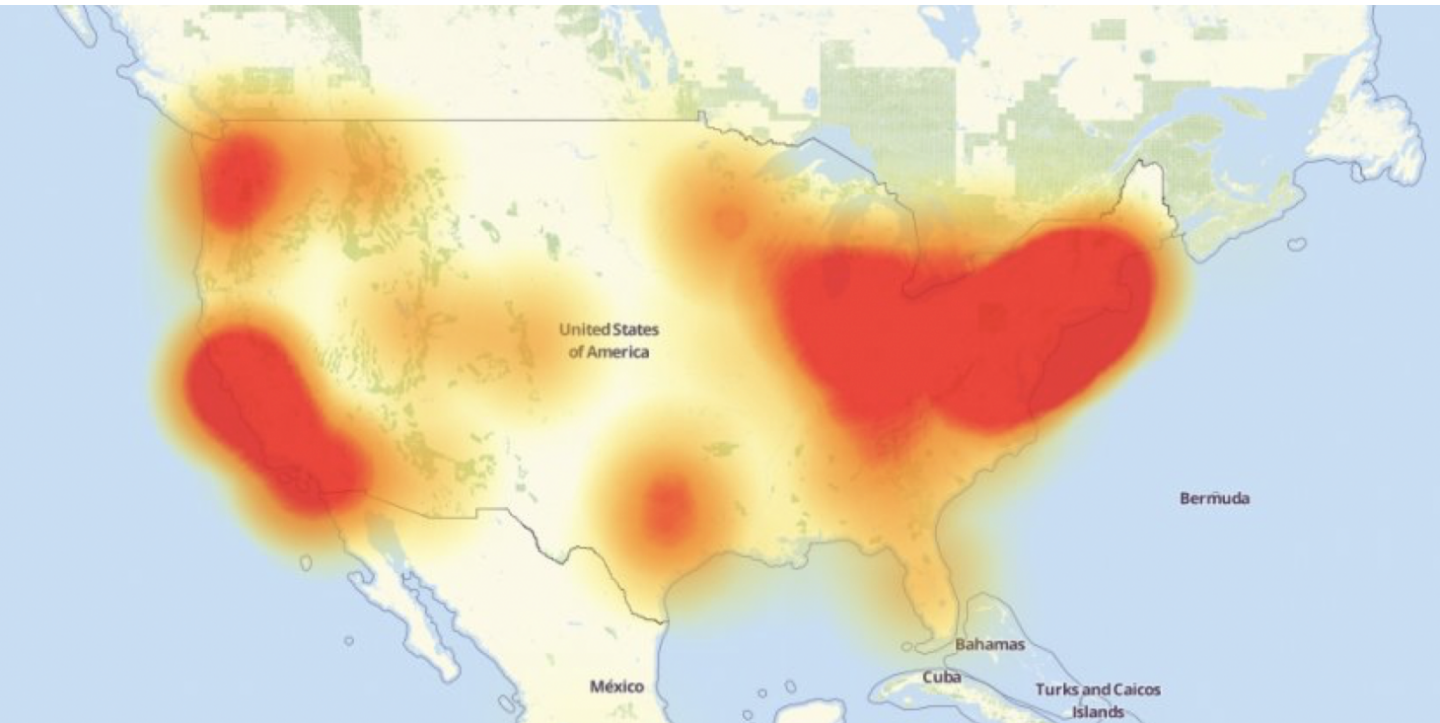
Due: 12/9

Presentations 12/5, 7, 9

Read Dataset Approval: 11/11?

Last few weeks of lab dedicated to project work.

A [DDoS attack](#) uses a variety of techniques to send countless junk requests to a website. This boosts traffic to the website so much that it gets overwhelmed, making it nearly impossible for anyone to load the page.



# ASSEMBLY – DE BRUIJN

Hamiltonian Path Problem

Eulerian Path Problem

# ASSEMBLY – DE BRUIJN

## Hamiltonian Path Problem

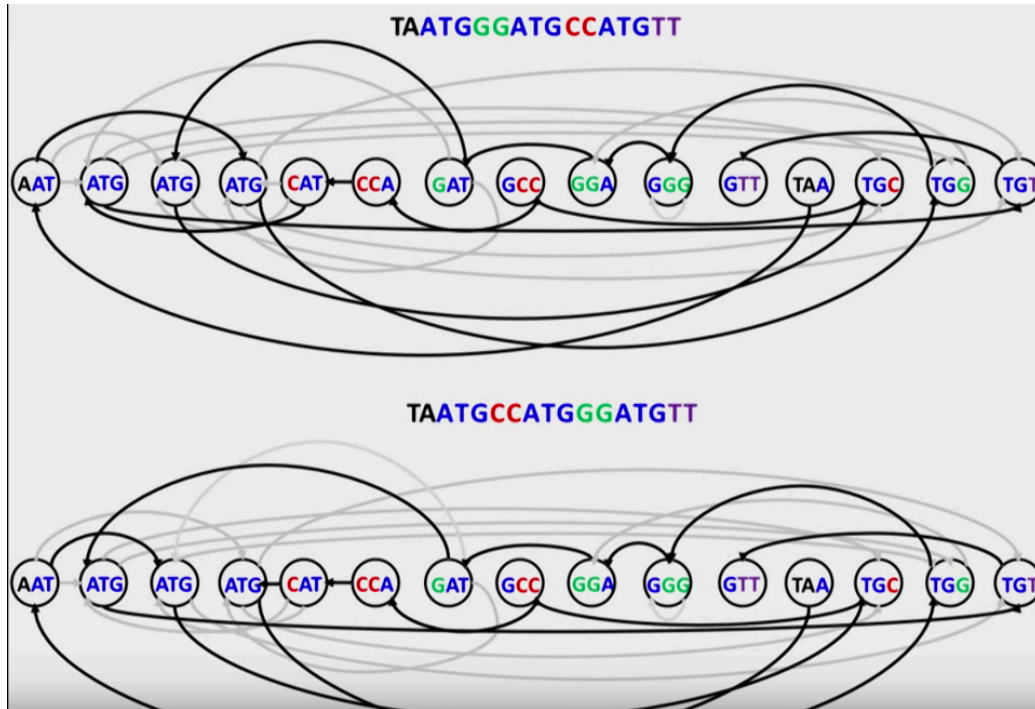
$Composition_3(TAATGCCATGGATGTT) =$



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

# ASSEMBLY – DE BRUIJN

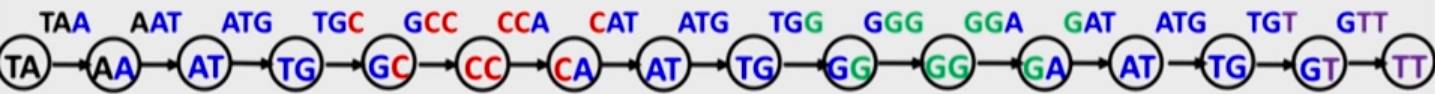
## Hamiltonian Path Problem





# ASSEMBLY – DE BRUIJN

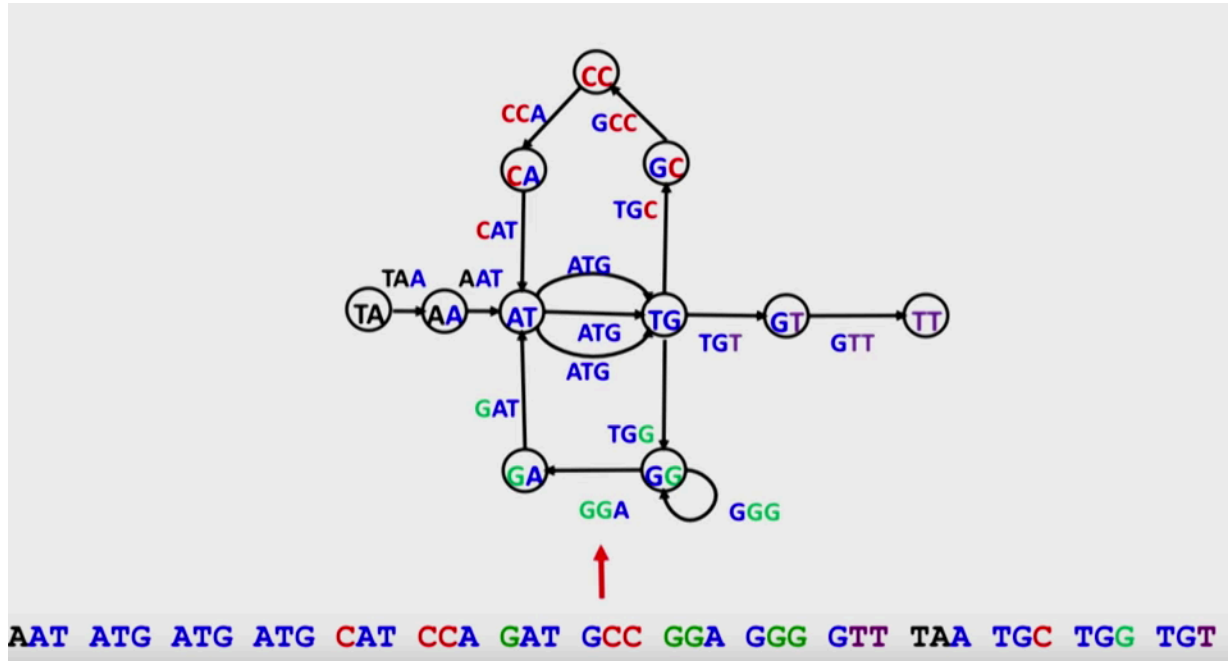
## Eulerian Path Problem



3-mers as **edges** and 2-mers as **nodes**

# ASSEMBLY – DE BRUIJN

## Eulerian Path Problem

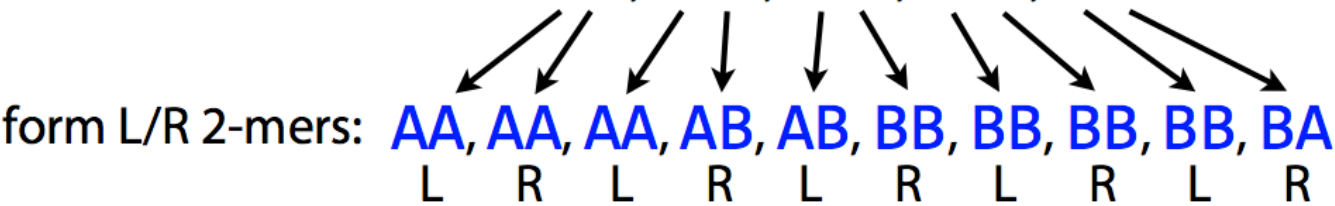


# ASSEMBLY – DE BRUIJN

# ASSEMBLY – DE BRUIJN

AAABBBBA

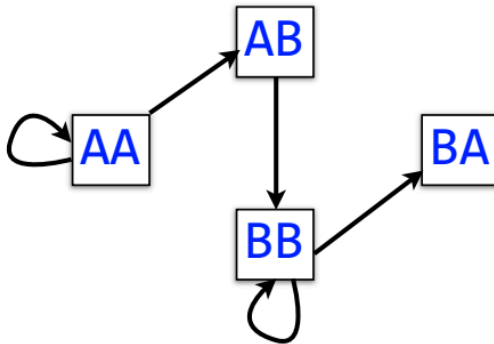
take all 3-mers: AAA, AAB, ABB, BBB, BBA



# 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

# ASSEMBLY – DE BRUIJN

GATTACAGTTCA

# ASSEMBLY – DE BRUIJN

GATTACAGTTCA  
GATTAC  
ACAGTTCA

# ASSEMBLY – DE BRUIJN

GATTAC

ACAGTTCA



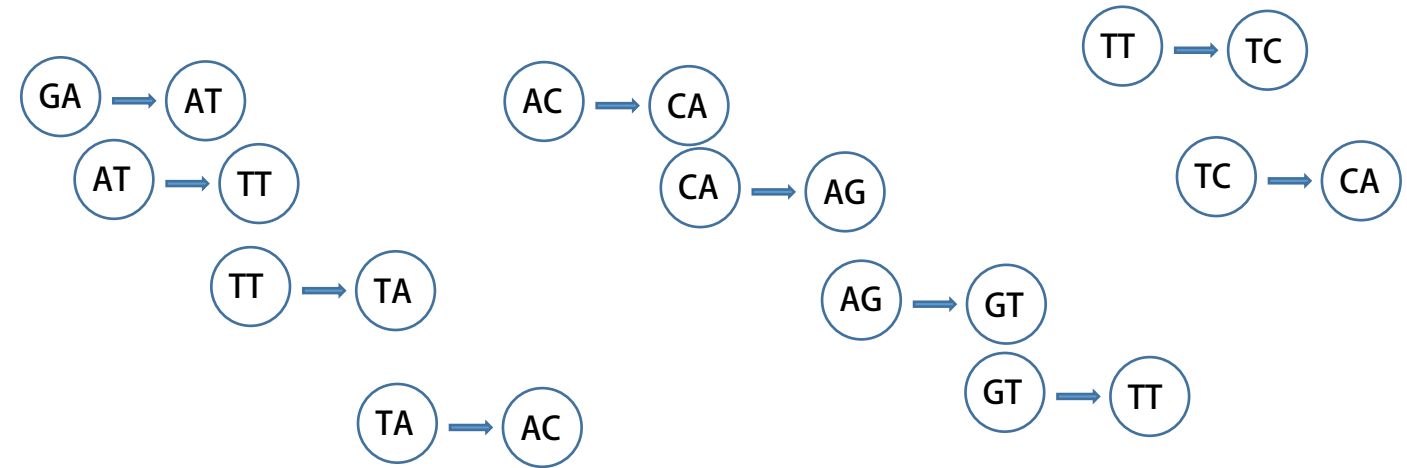
# ASSEMBLY – DE BRUIJN

GATTAC  
GAT  
ATT  
TTA  
TAC

ACAGTTCA  
ACA  
CAG  
AGT  
GTT  
TTC  
TCA

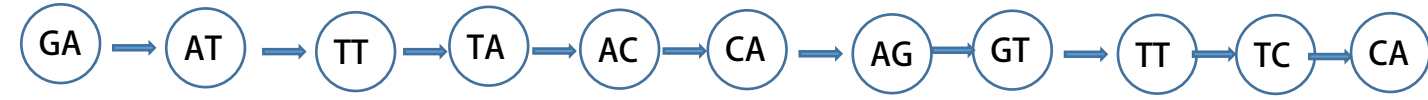
# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



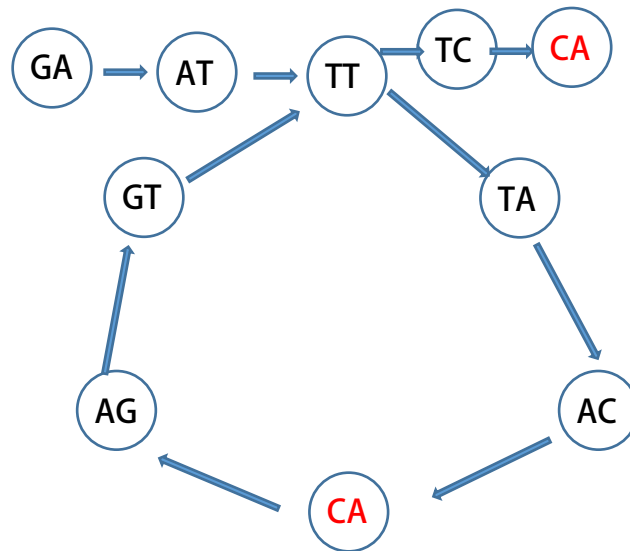
# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



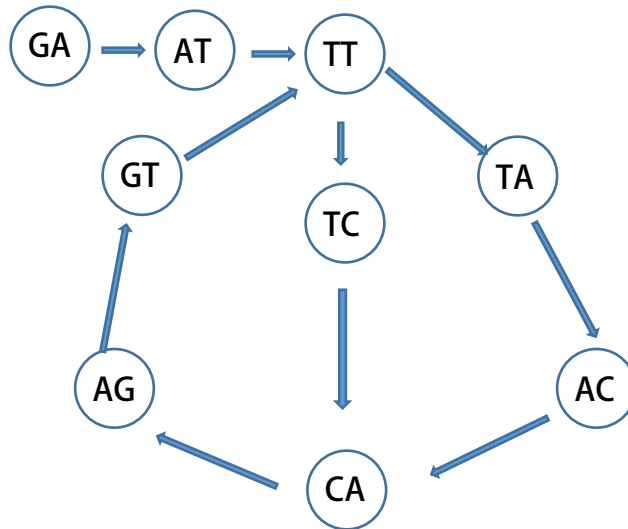
# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA

