

# **GENOME ASSEMBLY**

**26 OCT 15**

# ANNOUNCEMENTS

# ASSEMBLY – DE BRUIJN

## Hamiltonian Path Problem

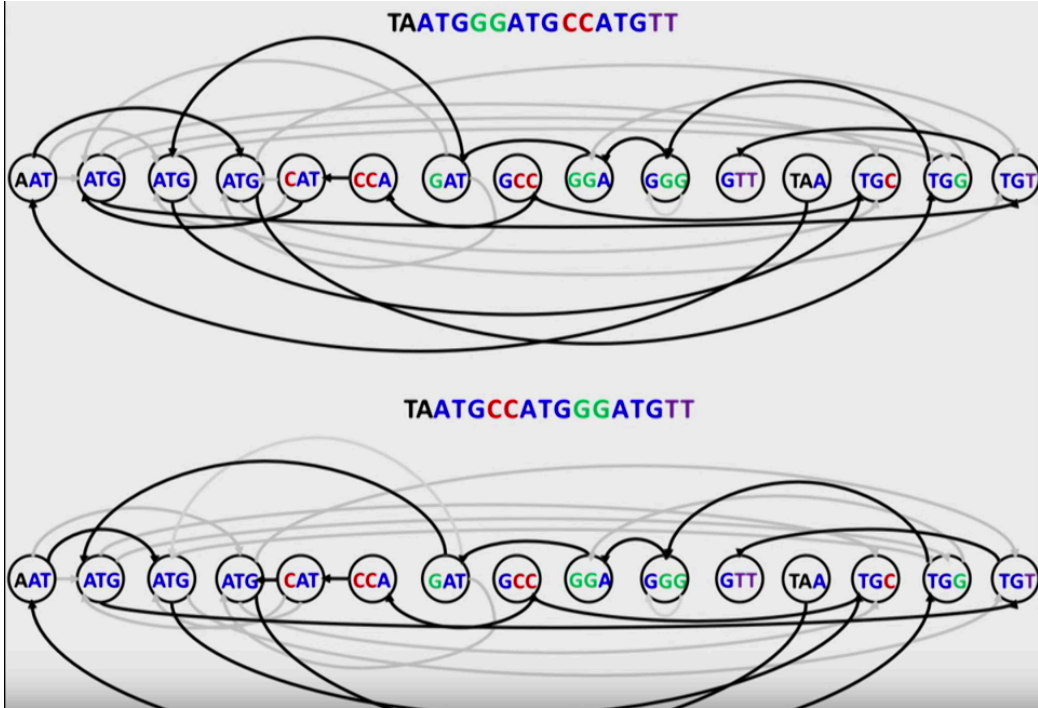
$Composition_3(TAATGCCATGGATGTT) =$



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

# ASSEMBLY – DE BRUIJN

## Hamiltonian Path Problem



# ASSEMBLY – DE BRUIJN

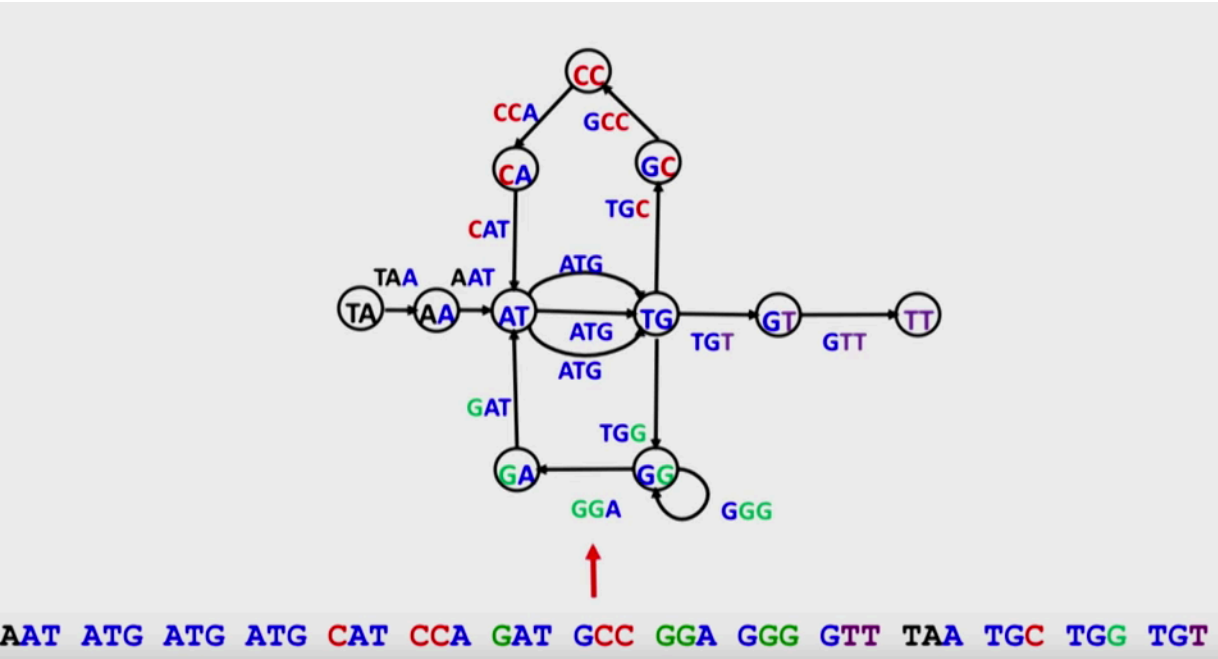
## Eulerian Path Problem



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

# ASSEMBLY – DE BRUIJN

## Eulerian Path Problem



# ASSEMBLY – DE BRUIJN

Hamiltonian vs Eulerian

# ASSEMBLY – DE BRUIJN GRAPH APPROACH



# ASSEMBLY – DE BRUIJN

AAABBBBA

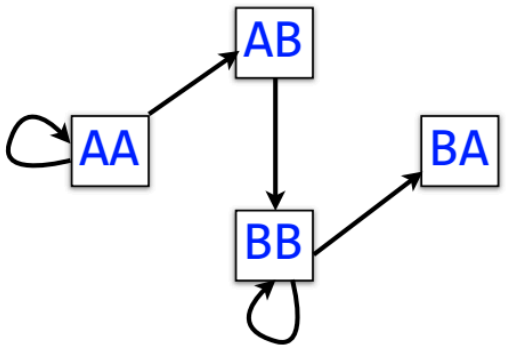
take all 3-mers: AAA, AAB, ABB, BBB, BBA

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

# 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

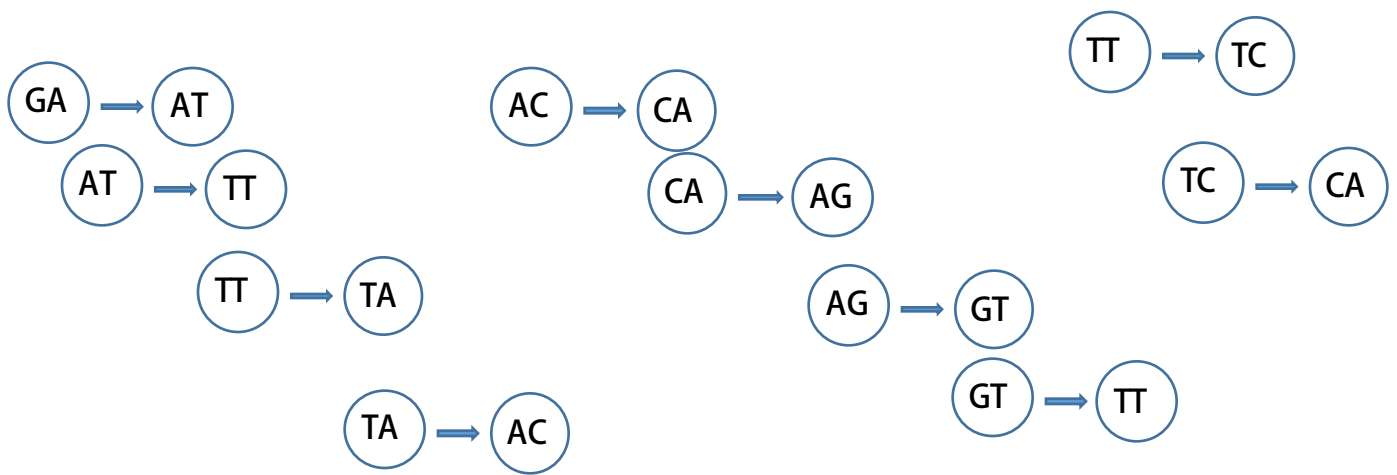
# ASSEMBLY – DE BRUIJN

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





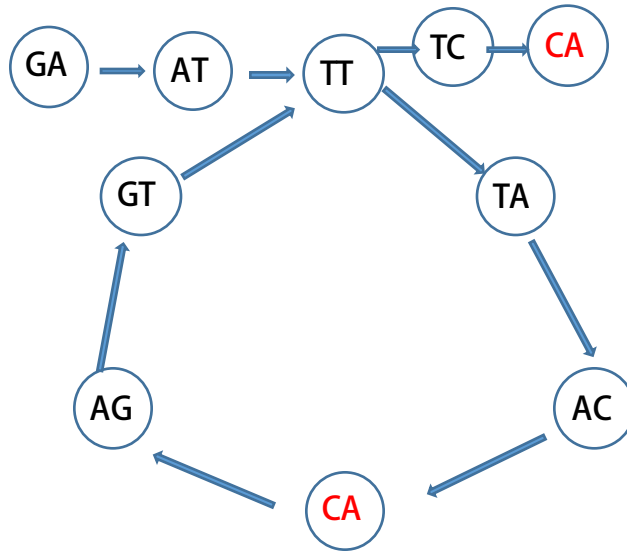
# ASSEMBLY – DE BRUIJN

GAT ATT TTA TAC ACA CAG AGT GTT TTC TCA



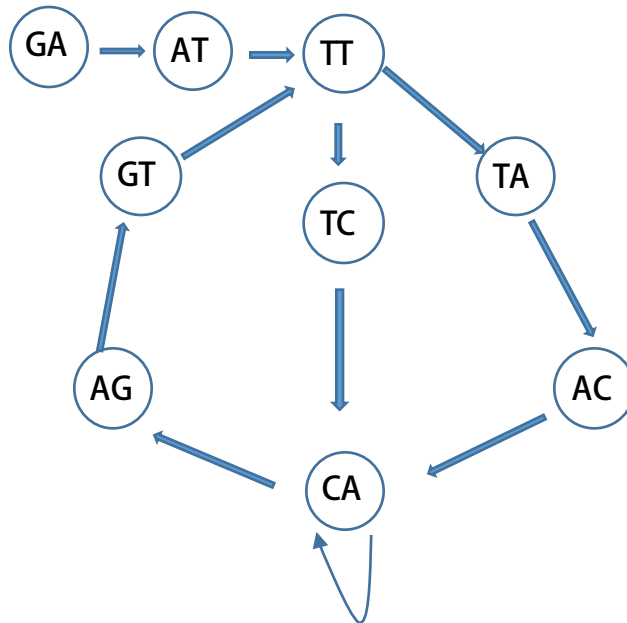
# ASSEMBLY – DE BRUIJN

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

