

Week 2 - 1/28/22

Friday, 21 January 2022 7:55 AM

$$Q = -10 \log_{10} P_e \quad ; \quad P_e = \text{P(error in base calling)}$$

eg. $Q = 30$; ASCII (?)

$$30 = -10 \log_{10} P_e$$

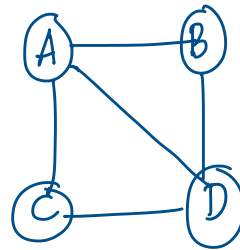
$$-3 = \log_{10} P_e$$

$$\Rightarrow 10^{-3} = P_e = 0.001$$

GRAPH THEORY 'REVIEW'

OBJECT - NODE

RELATIONSHIP - EDGES



$$\text{NODES} = \{A, B, C, D\}$$

$$\text{EDGES} = \{AB, BD, CD, AC, AD\}$$

Problem 1 : find a path connecting all nodes ; such that every node is visited ONLY once [HAMILTONIAN PATH]

① $A \rightarrow B \rightarrow D \rightarrow C$

② $A \rightarrow C \rightarrow D \rightarrow B$

HAMILTONIAN CYCLE :

$A \rightarrow B \rightarrow D \rightarrow C \rightarrow A$

Problem 2: find a path such that every edge is visited ONLY ONCE [EULERIAN PATH]

$A \rightarrow B \rightarrow D \rightarrow C \rightarrow A \rightarrow D$

GENOME : ATGGCGTGCA ✓
(CIRCULAR)

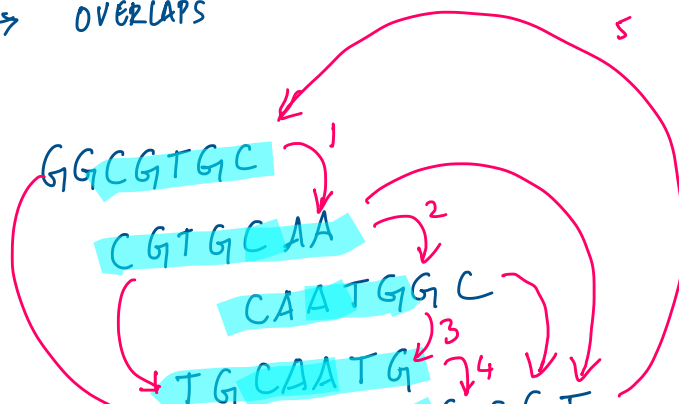
READS : CGTGCAA, ATGGCGT, CAATGGC,

GGCGTGC, TGCAATG

consider each read as a NODE

edges \rightarrow 'OVERLAPS'

HAMILTONIAN CYCLE





 ATG G C G T

 G G C G T G C A A T

'k-mer length' \rightarrow length of overlap between reads

de Bruijn graphs ; Eulerian path

SUPERSTRING PROBLEM

\hookrightarrow find the shortest superstring that contains all possible substrings of length = k [k-mer] in a given alphabet.

eg. ALPHABET = $\{0, 1\}$; $n = 2$ [# of letters]

k-mer length = 3

$\{000, 001, 010, 100, 110, 101, 011, 111\}$

n^k possible k-mers

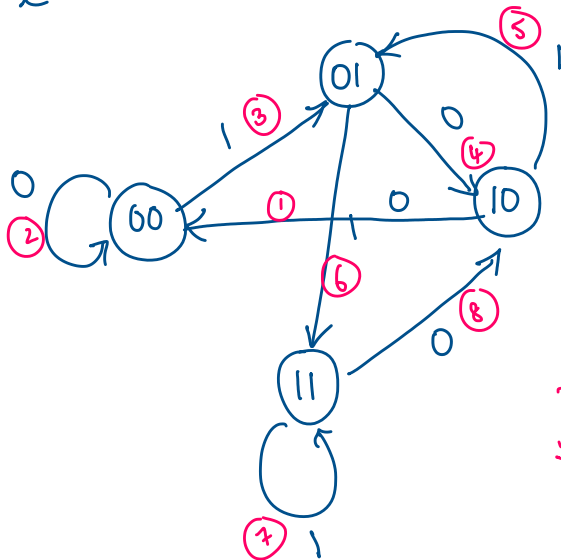
$$\left[\begin{matrix} 3 \\ 2 \end{matrix} = 8 \right]$$

de Bruijn graph

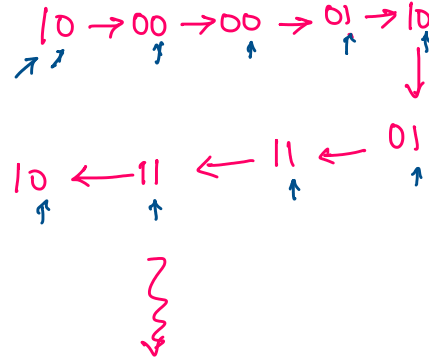
Construct a de Bruijn graph

such that

each node = ' $k-1$ ' mer



EULERIAN
PATH



100010110

ALPHABET = $\{A, C, G, T\}$ $n=4$ $k\text{-mer} = 3$

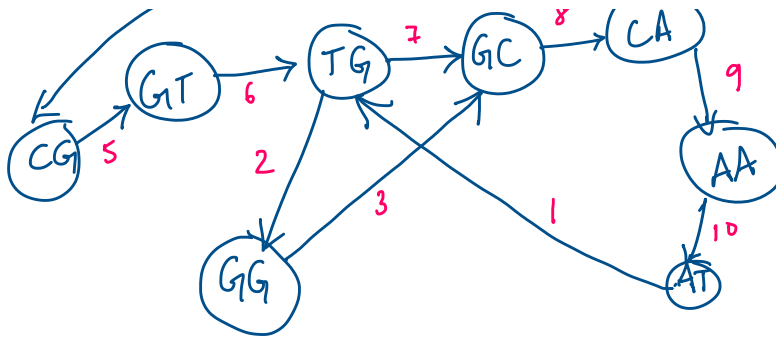
de Bruijn graph has to have nodes of
length = $k-1 = 3-1 = 2$

possible nodes = $\{AA, AC, AG, AT, \dots, TT\}$
[4]

OBSERVED

READS: $\left[\begin{array}{l} \underline{CGTGC}AA, \underline{TGCAATG}, \underline{ATGGCGT}, \\ \underline{GGCGTGC}, \underline{CAATGGC} \end{array} \right]$

nodes = $\{ \underline{CG}, \underline{GT}, \underline{TG}, \underline{GC}, \underline{CA}, \underline{AA}, \underline{AT}, \underline{GG} \}$



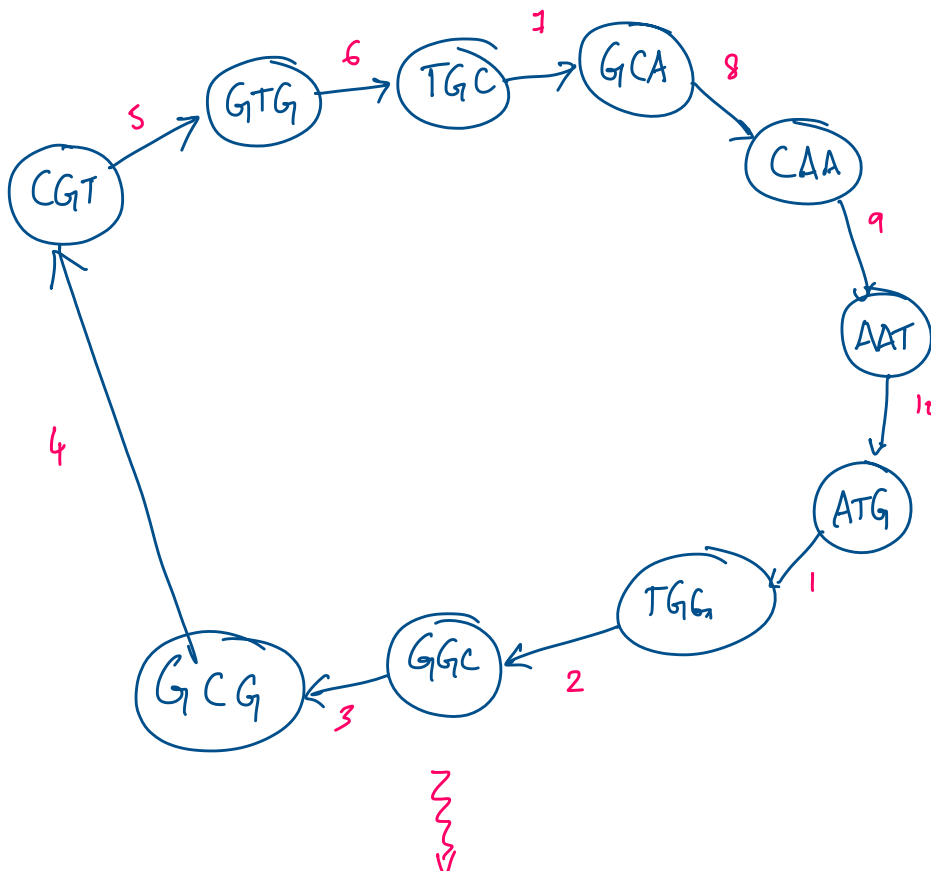
EULERIAN CYCLE:



ASSEMBLED
GENOME

$k = 4$

nodes : $\{ \text{CGT, GTG, TGC, GCA, CAA, ATG, TGC, GGC, GCG} \}$
AAT,



ATGGCGTGCA

