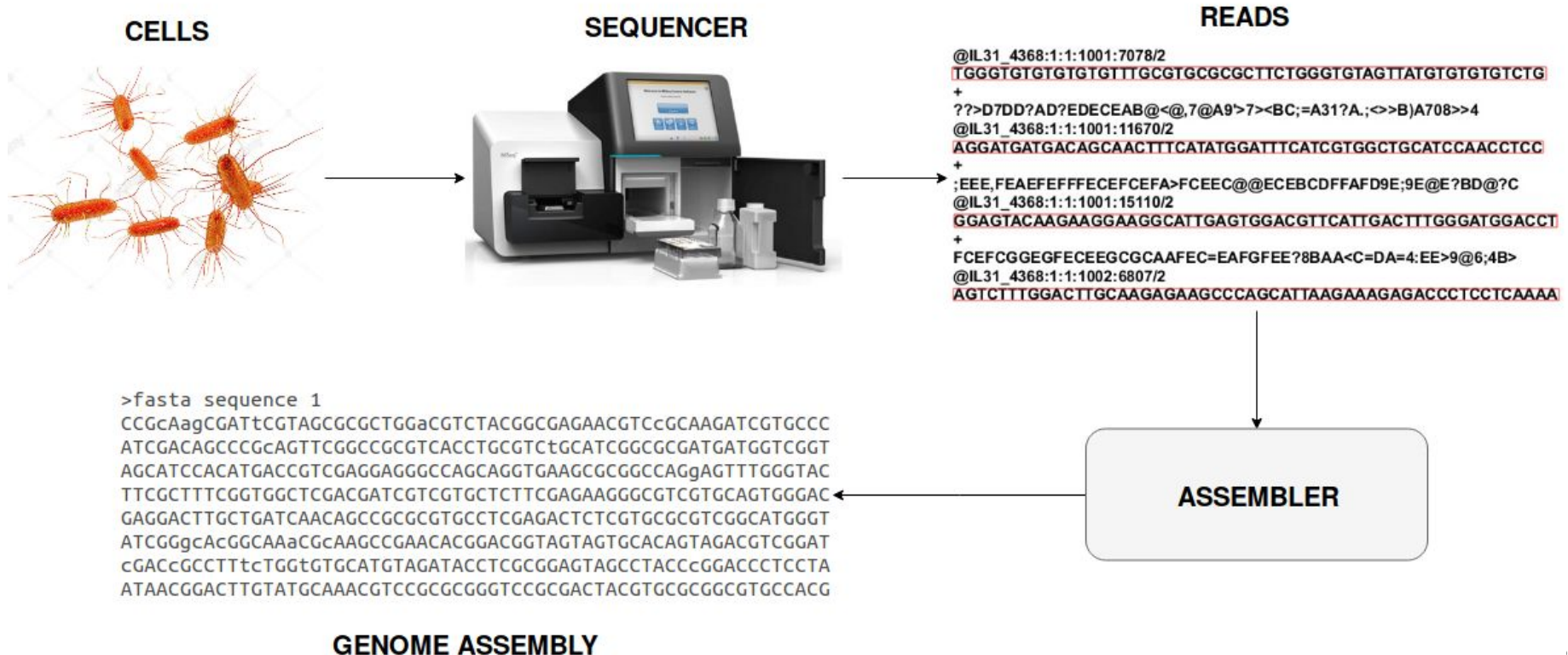


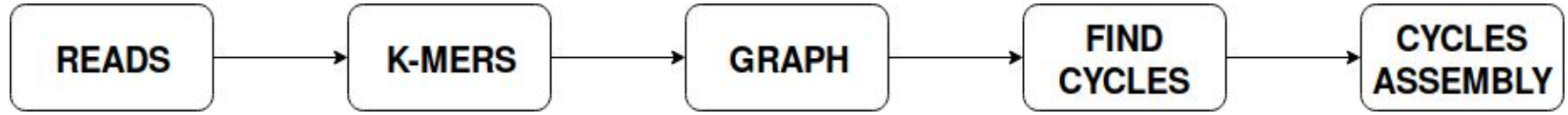
# De Novo Genome Assembly

Etienne JEAN  
M2 BI - Short Project Presentation  
September 2018

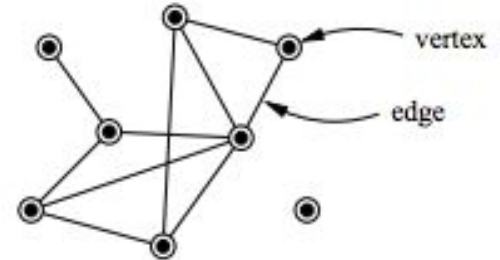
# Principle of de novo assembly



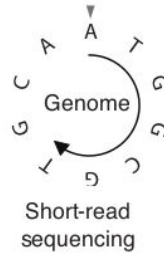
# Description of an assembler



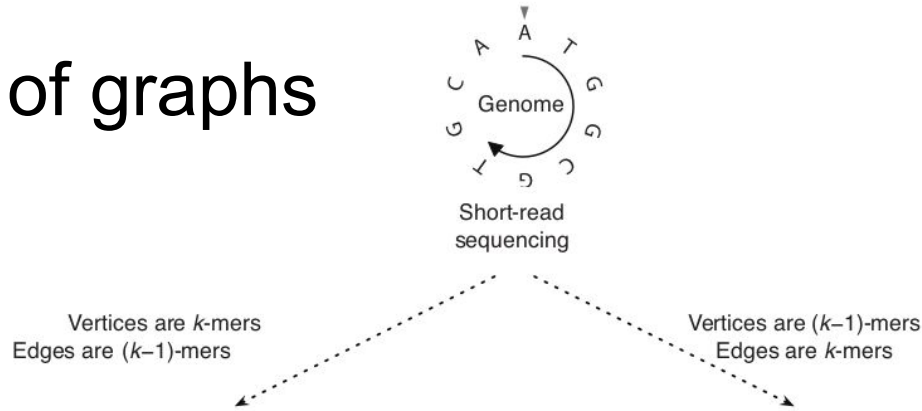
- **K-mers** strategy : substrings of length K
  - Remove redundancy of reads
  - Graph is no bigger than needed
  - Necessary step for De Bruijn graph
- **Graph** : set of **vertices** connected by **edges**



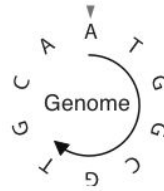
# Two types of graphs



# Two types of graphs

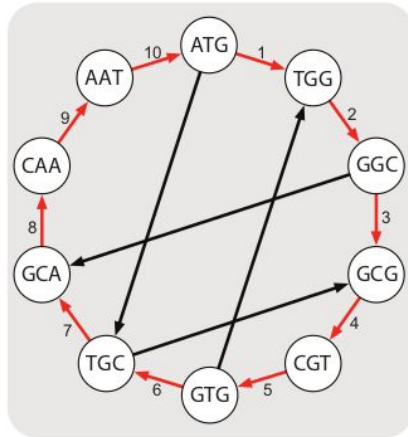


# Two types of graphs



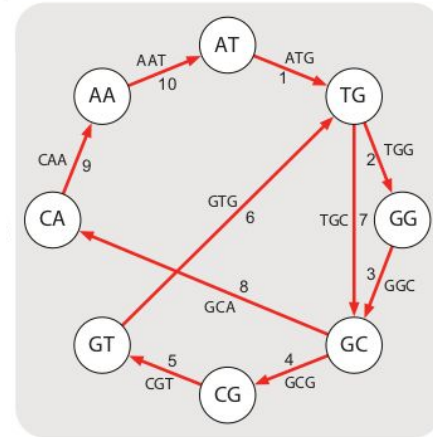
Short-read  
sequencing

Vertices are  $k$ -mers  
Edges are  $(k-1)$ -mers



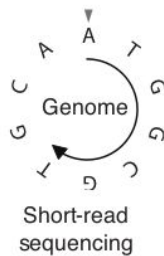
**Hamiltonian cycle**  
Visit each vertex once  
(harder to solve)

Vertices are  $(k-1)$ -mers  
Edges are  $k$ -mers

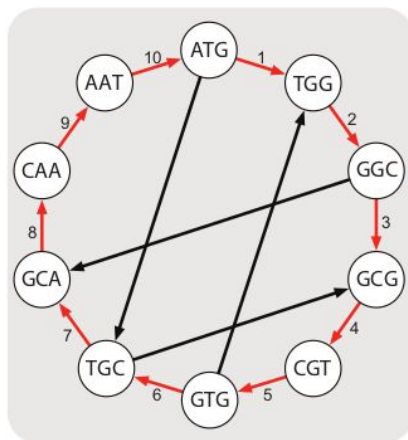


**Eulerian cycle**  
Visit each edge once  
(easier to solve)

# Two types of graphs

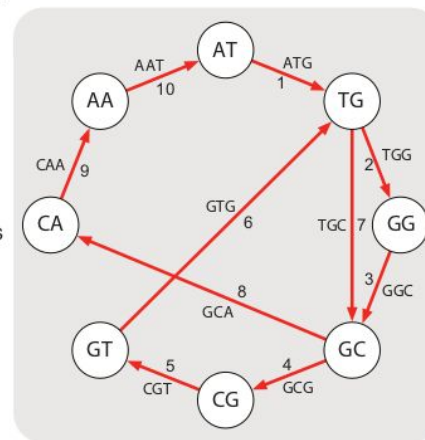


Vertices are  $k$ -mers  
Edges are  $(k-1)$ -mers



**Hamiltonian cycle**  
Visit each vertex once  
(harder to solve)

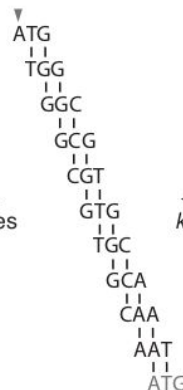
Vertices are  $(k-1)$ -mers  
Edges are  $k$ -mers



**Eulerian cycle**  
Visit each edge once  
(easier to solve)

$k$ -mers from vertices

$k$ -mers from edges



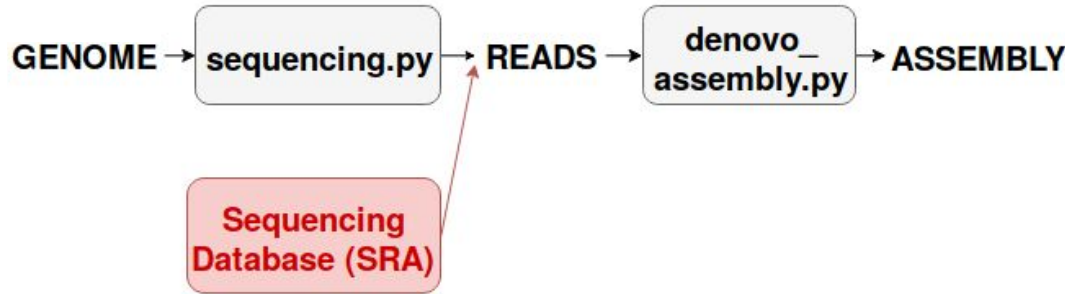
Genome: ATGGCGTGCAATG

# Program architecture : 4 different modules

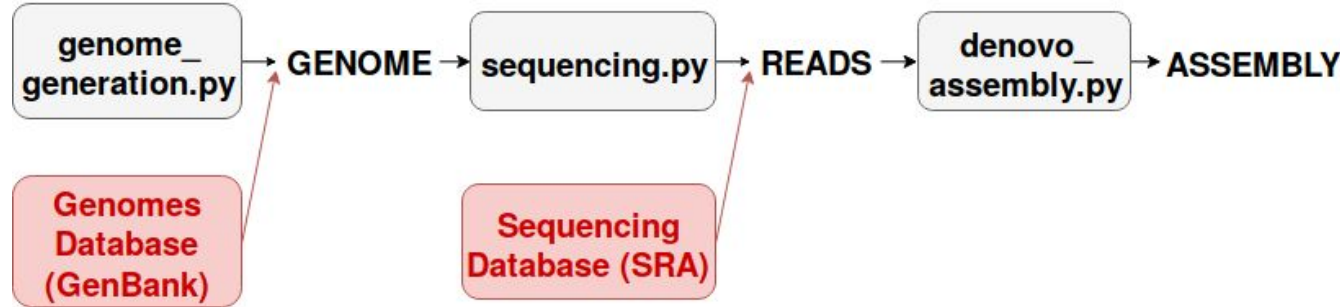




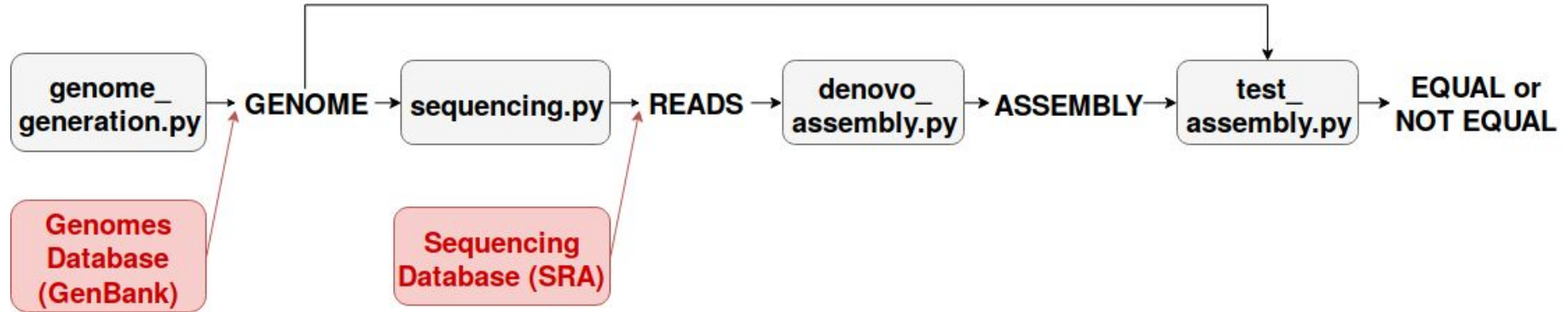
# Program architecture : 4 different modules



# Program architecture : 4 different modules

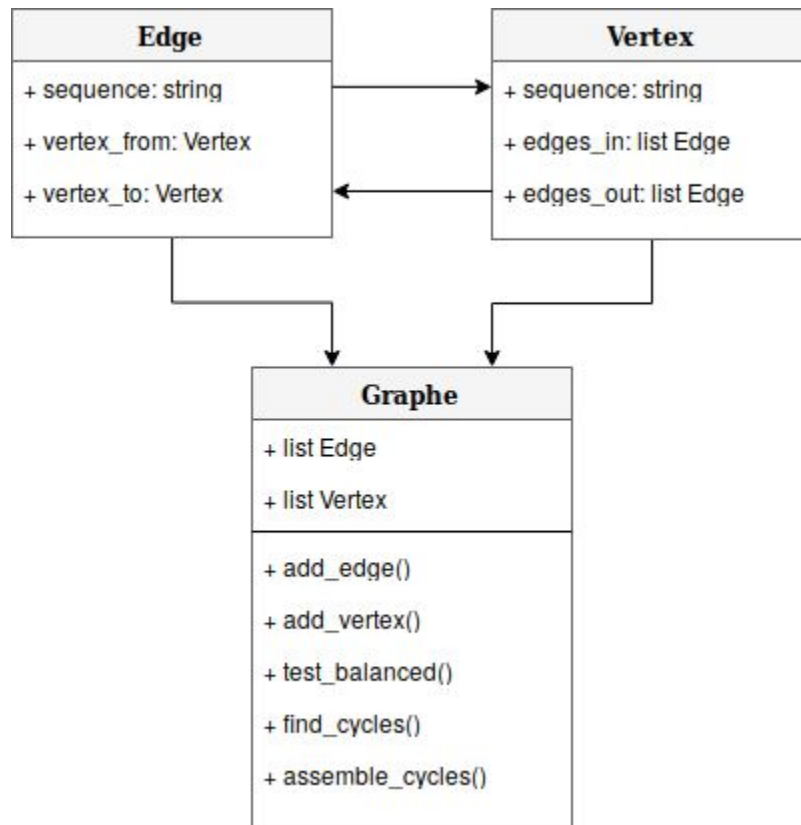


# Program architecture : 4 different modules



# Graph implementation

- Object oriented
- Vertices object point to Edges object
  - Reduces computation time
- Adjacency list
  - More adapted to sparse graphs than adjacency matrix
- Test for balanced De Bruijn graph
- Find Eulerian cycles
- Assemble cycles into one



# Results

- Random genome, 1 Mb :

- reads 100 bp
- coverage 50x
- K-mers 55 bp

93 seconds

- Mycoplasma genitalium genome, 586 kb :

- reads 400 bp
- coverage 100x
- K-mers 250 bp

82 seconds

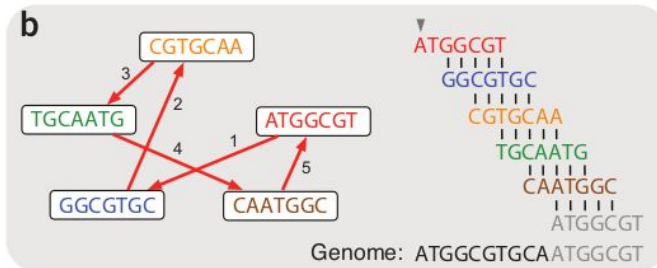
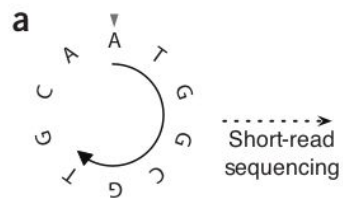
# Improvements

- Sequencing errors management
  - sequencing.py is currently not optional
- Multiple assemblies
  - 2 cycles connected by 2 different nodes
- Linear chromosome assembly
  - Eulerian path instead of Eulerian cycle
- Partial assemblies
  - When graph is not entirely balanced
  - Find set of contigs (subsequences of the genome)
- Multiple linear chromosomes assembly
  - Find different contigs, that correspond to different chromosomes
- Repeated sequences management
  - Introduce k-mers multiplicity

Thank you for your attention

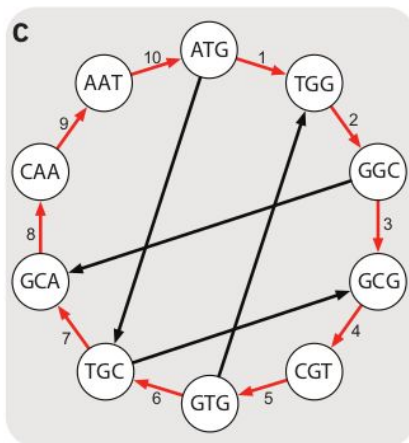
# Supplementary slides



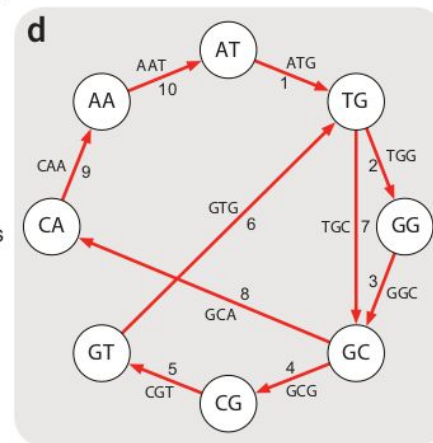
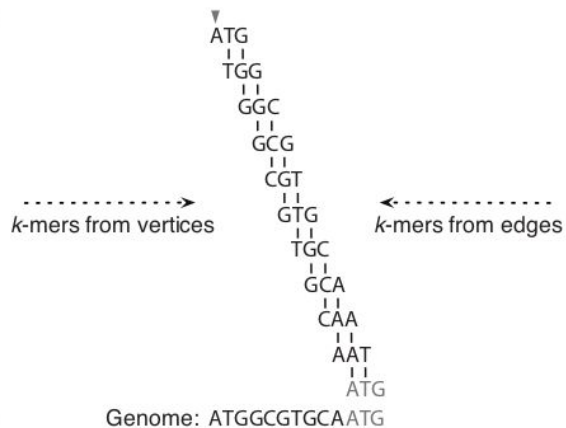


Vertices are  $k$ -mers  
Edges are pairwise alignments

Vertices are  $(k-1)$ -mers  
Edges are  $k$ -mers

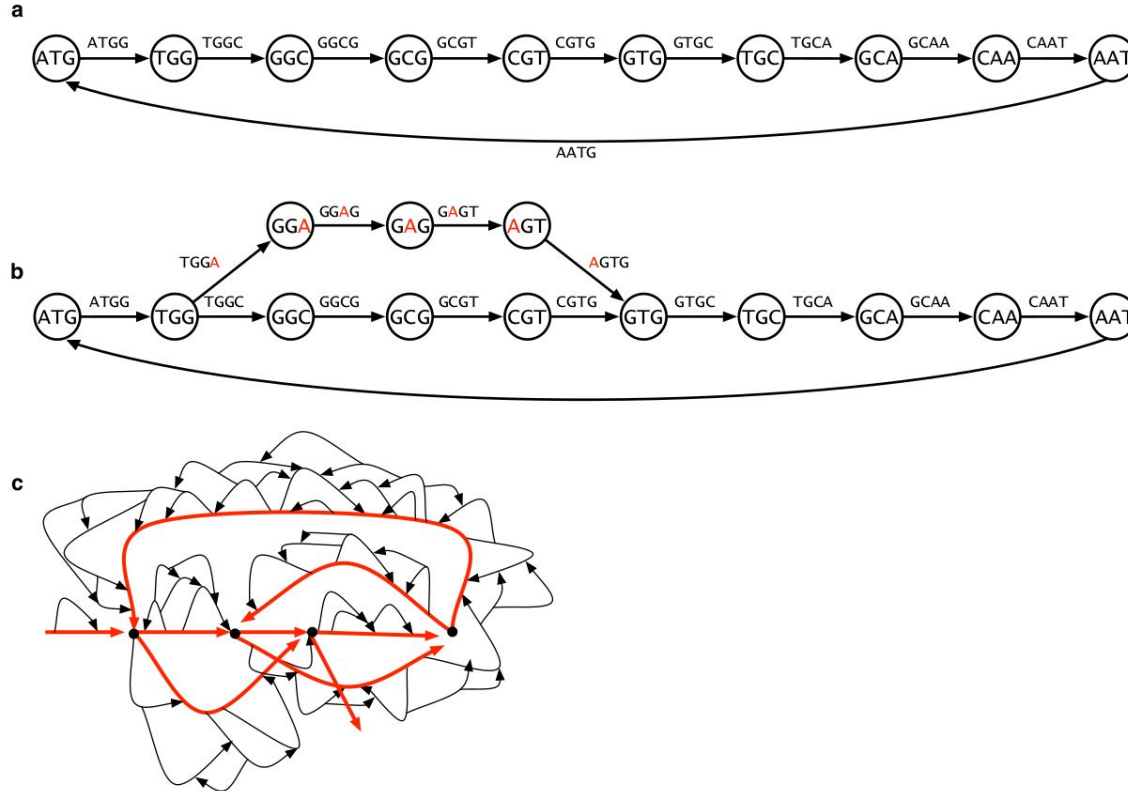


**Hamiltonian cycle**  
Visit each vertex once  
(harder to solve)



**Eulerian cycle**  
Visit each edge once  
(easier to solve)

# Sequencing errors : creation of bulges in the graph



# K-mers multiplicity for repeated sequences

