

Kako složiti genomske slagalice od miliona delova?

Grafovski algoritmi

*Bioinformatics Algorithms:
an Active Learning Approach*
~Poglavlje 3~

Pregled

- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- Problem rekonstrukcije niske
- Rekonstrukcija niske kao problem Hamiltonove putanje
- Rekonstrukcija niske kao problem Ojlerove putanje
- De Bruijinovi grafovi
- Ojlerova teorema
- Spajanje parova očitavanja
- U realnosti

Genom

- Genom jednog organizma predstavlja njegov genetski materijal
- Kod većine organizama, genetski materijal je sadržan u DNK
- Kod čoveka, genom sadrži oko tri milijarde nukleotida
- Genomi nekih organizama su i 100 puta veći od humanog genoma

Amoeba Dubia
~ 670 milijardi



Jovana Kovačević, Bioinformatika



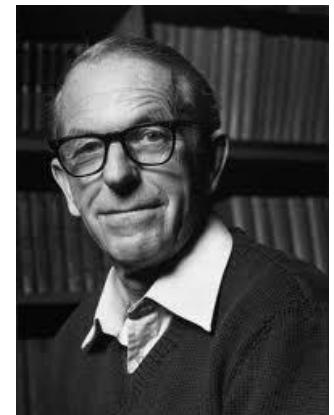
Paris Japonica
~ 150 milijardi

Kratka istorija sekpcioniranja genoma

- **1977:** Walter Gilbert i Frederick Sanger razvijaju nezavisne metode sa sekpcioniranje DNK
- **1980:** Podelili su Nobelovu nagradu.
- Njihove metode za sekpcioniranje su bile veoma skupe (\$3 milijarde za sekpcioniranje humanog genoma).

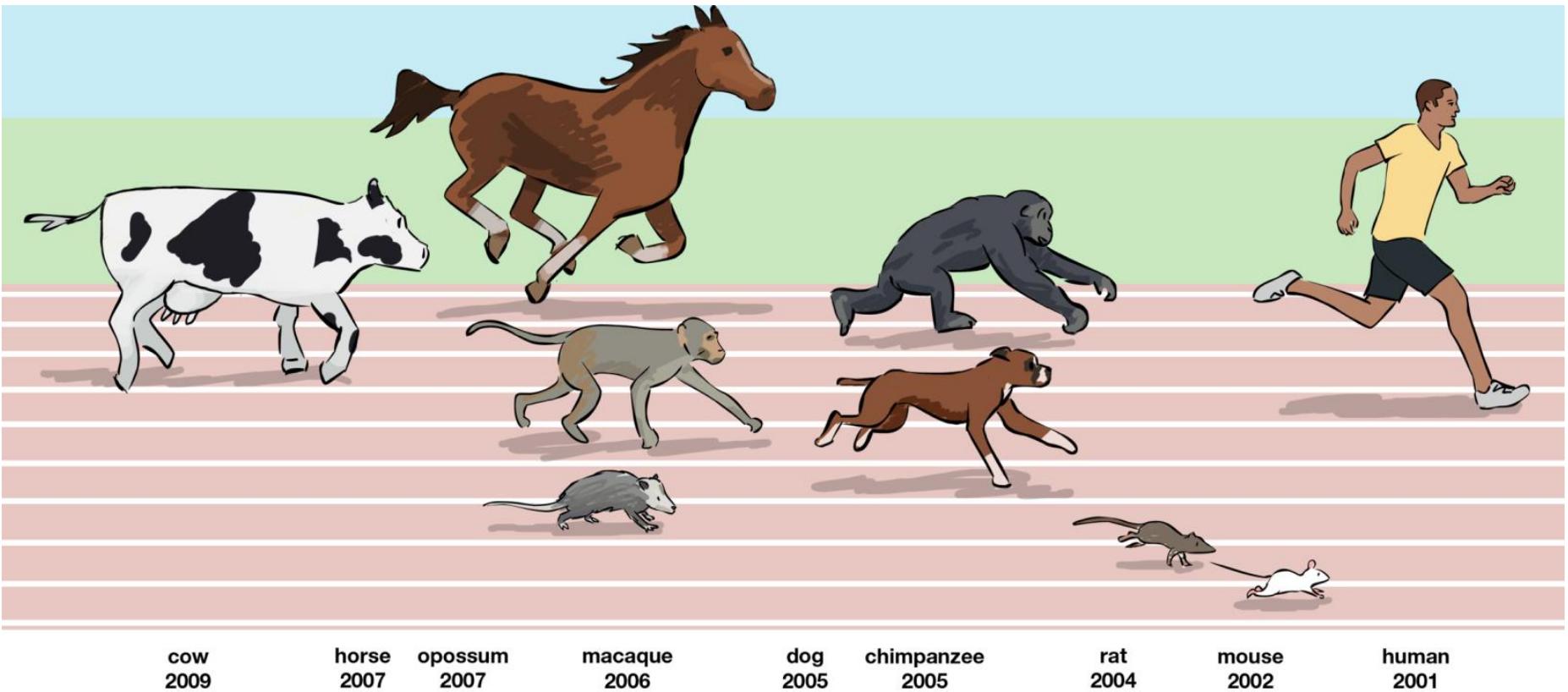


Walter Gilbert



Frederick Sanger

Početak sekvencioniranje genoma



- Krajem 2000-tih Sanger metodom je sekvencioniran veliki broj genoma
- Visoka cena je bila ograničavajući faktor i za dalji napredak je bila neophodna nova tehnologija sekvencioniranja

Sekvencioniranje nove generacije

- *Next Generation Sequencing (NGS)*
- **Krajem 2000-tih:** Na tržištu se pojavljuju nove mašine za sekvencioniranje
 - *Illumina* smanjuje trošak sekvencioniranja humanog genoma sa 3 milijarde na 10 hiljada dolara
 - Kompanija *Complete Genomics* otvara genomsku fabriku u Silikonskoj dolini koja sekvencionira stotine genoma mesečno
 - Pekinški genomski institut (*BGI - Beijing Genome Institute*) preuzima *Complete Genomics* 2013. godine i postaje najveći svetski centar za sekvencioniranje genoma



illumina

Complete
genomics

华大基因
BGI

Sekvencioniranje ličnih genoma



Sekvencioniranje ličnih genoma

- Genomi se kod različitih ljudi razlikuju na malom broju pozicija (u proseku sadrže jednu mutaciju na hiljadu nukleotida)
- Ova razlika je odgovorna za različite visine kod ljudi, da li će imati sklonost ka visokom holesterolu ili ne, za veliki broj genetskih bolesti, ...



```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGA  
TCAGCTACCACATCGTAGCTACGATGCATTAGCAAGCTAT  
CGATCGATCGATCGATTATCTACGATCGATCGATCGATCA  
CTATACGAGCTACTACGTACGTACGATCGCGGGACTATT  
TCGACTACAGATAAAACATGCTAGTACAACAGTATACATA  
GCTGCGGGATACGATTAGCTAATAGCTGACGATATCCGAT
```



```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGA  
TCAGCTACACACATCGTAGCTACGATGCATTAGCAAGCTAT  
CGATCGATCGATCGATTATCTACGATCGATCGATCGATCA  
CTATACGAGCTACTACGTACGTACGATCGCGTGACTATT  
TCGACTACAGATGAAACATGCTAGTACAACAGTATACATA  
GCTGCGGGATACGATTAGCTAATAGCTGACGATATCCGAT
```

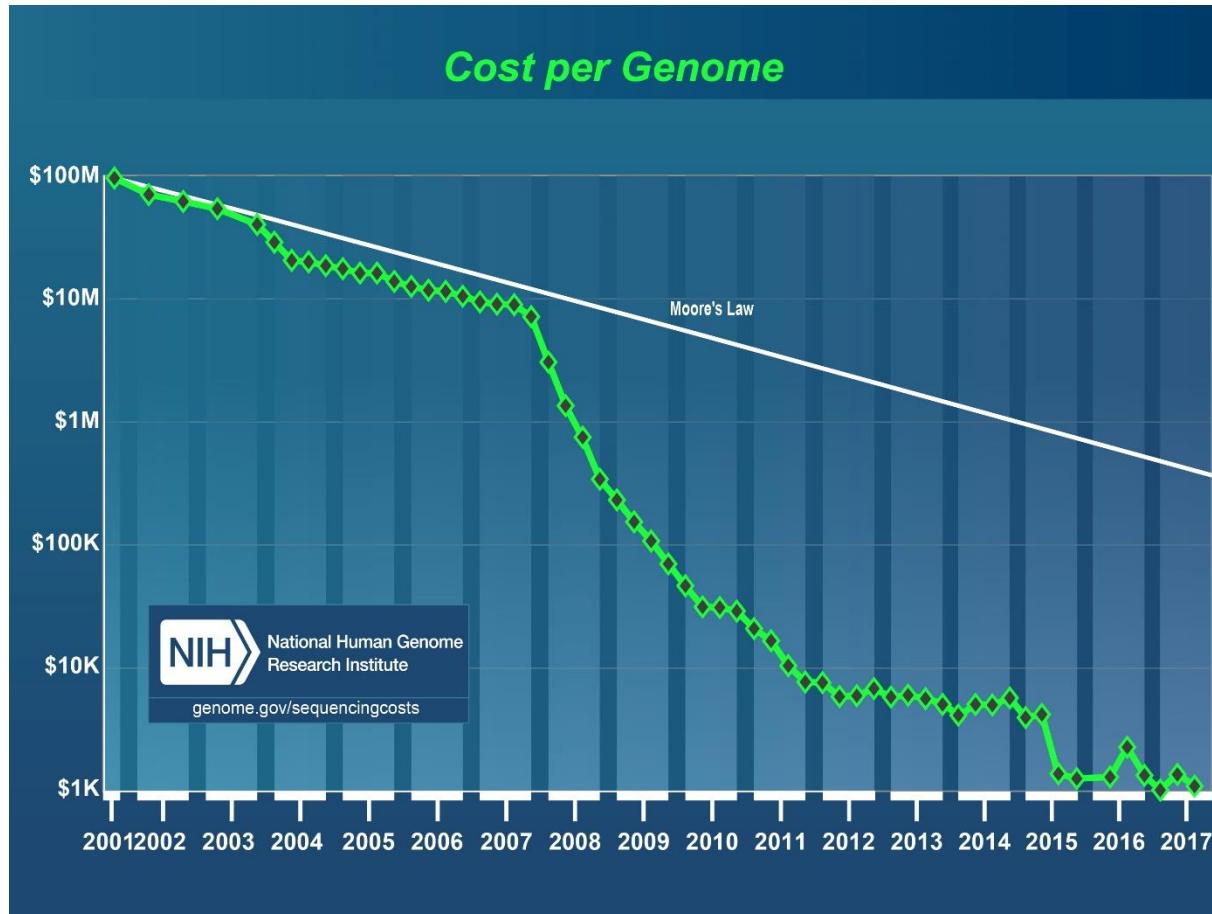


Sekvencioniranje ličnih genoma

- **2010:** *Nicholas Volker* je postao prvo ljudsko biće čiji je život spašen zahvaljujući genomskom sekvencioniranju
 - Lekari nisu mogli da postave tačnu dijagnozu i morali su da ga podvrgnu velikom broju operacija pokušavajući da je utvrde
 - Sekvencioniranje je otkrilo retku mutaciju na jednom genu (*XIAP*) koja je bila povezana sa oštećenjem njegovog imunog sistema
 - Ovo otkriće je navelo lekare na adekvatnu terapiju koja je rešila problem



Sekvencioniranje ličnih genoma



Pregled

- Šta je sekvencioniranje genoma?
- **Eksplozija u štampariji**
- Problem rekonstrukcije niske
- Rekonstrukcija niske kao problem Hamiltonove putanje
- Rekonstrukcija niske kao problem Ojlerove putanje
- De Bruijinovi grafovi
- Ojlerova teorema
- Spajanje parova očitavanja
- U realnosti

Problem novina

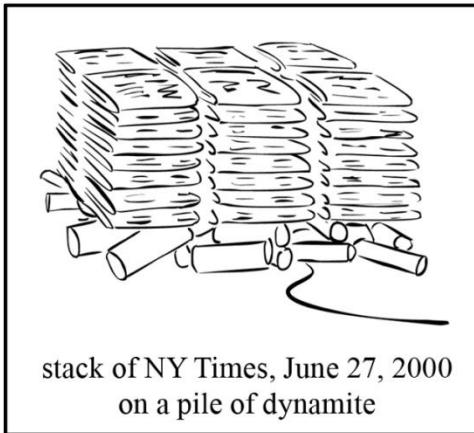


stack of NY Times, June 27, 2000

Problem novina

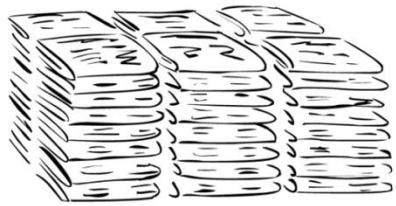


stack of NY Times, June 27, 2000



stack of NY Times, June 27, 2000
on a pile of dynamite

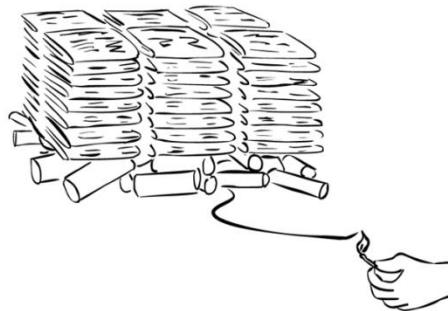
Problem novina



stack of NY Times, June 27, 2000

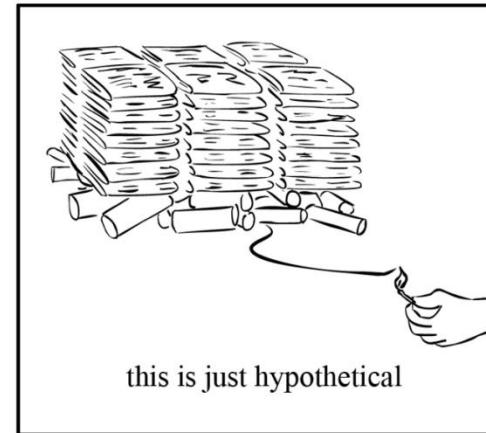
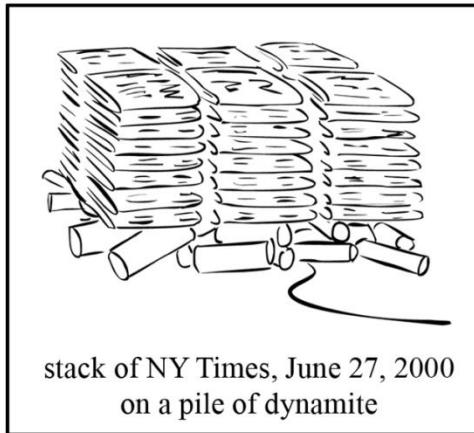


stack of NY Times, June 27, 2000
on a pile of dynamite



this is just hypothetical

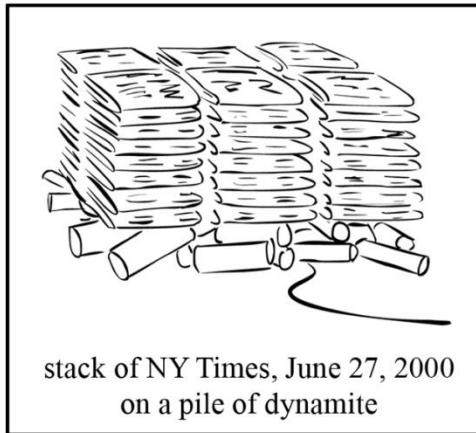
Problem novina



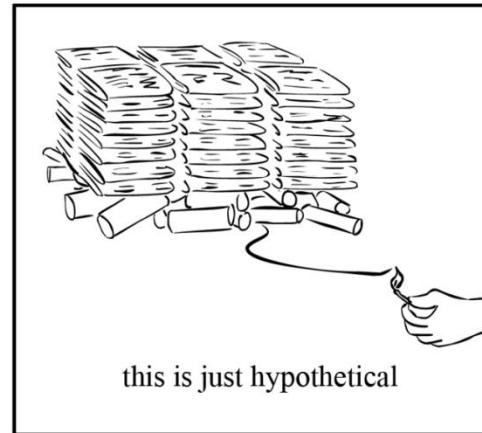
Problem novina



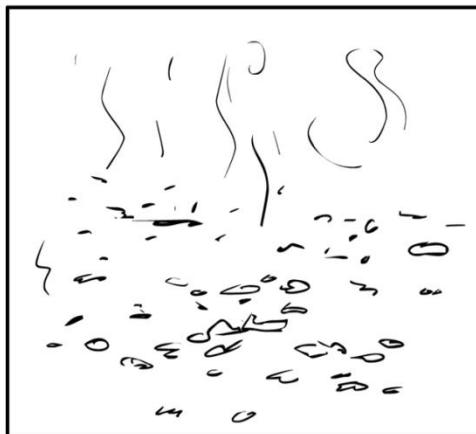
stack of NY Times, June 27, 2000



stack of NY Times, June 27, 2000
on a pile of dynamite



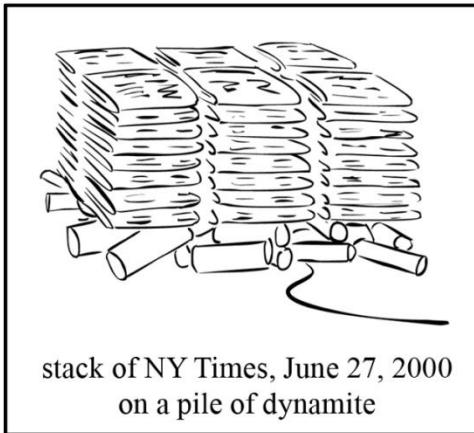
this is just hypothetical



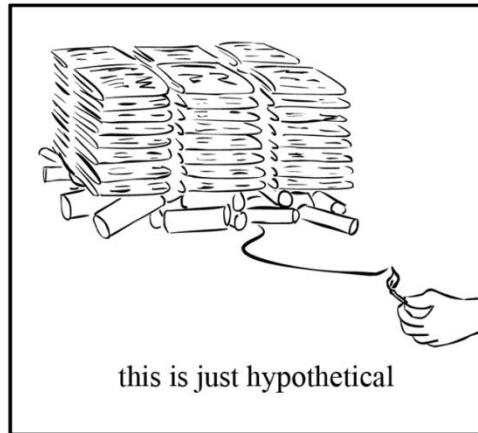
Problem novina



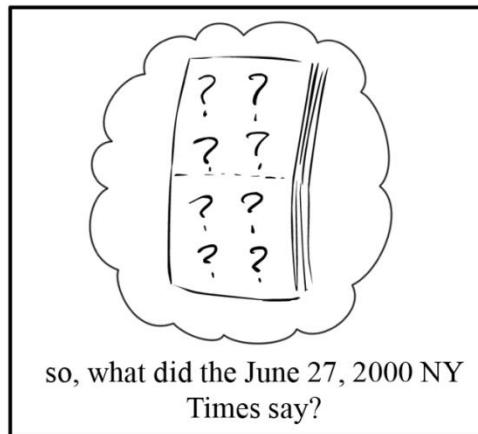
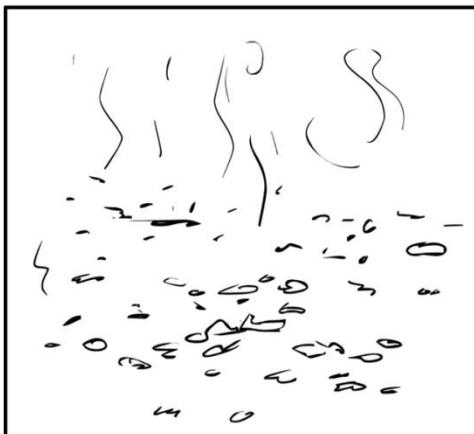
stack of NY Times, June 27, 2000



stack of NY Times, June 27, 2000
on a pile of dynamite

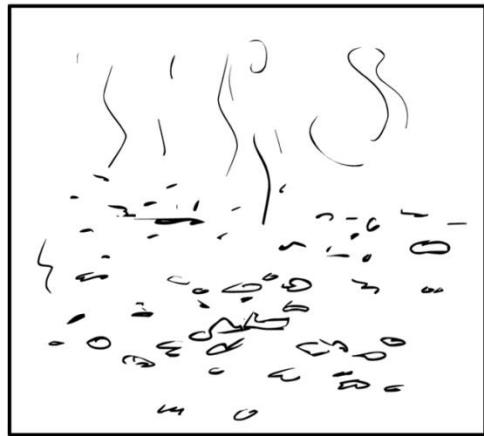


this is just hypothetical



so, what did the June 27, 2000 NY
Times say?

Problem novina kao delovi slagalice koji se preklapaju

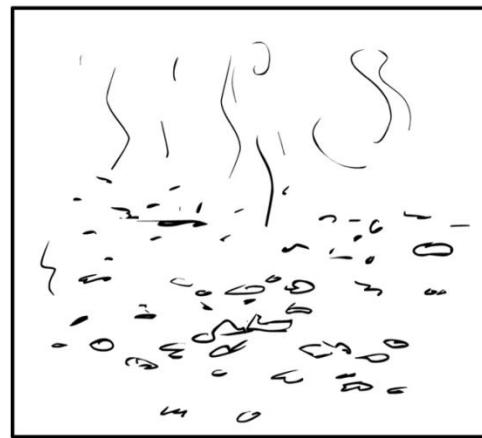


hoodie, appi
we have not yet named
information is welc

lie, appi
yet named any suspects, alt
is welc

o'2
ce ca

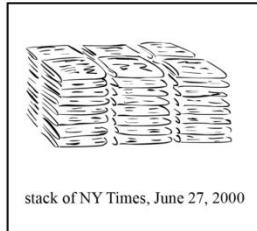
Problem novina kao delovi slagalice koji se preklapaju



hoodie, appre
ce have not yet named
information is welc

die, appre
yet named any suspects, alt
is welc

Milion kopija genoma



stack of NY Times, June 27, 2000

CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTACGATGCATTAGCAAGCTATCGGATCAGCTACCACATCGTAGC
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTACGATGCATTAGCAAGCTATCGGATCAGCTACCACATCGTAGC
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTACGATGCATTAGCAAGCTATCGGATCAGCTACCACATCGTAGC
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACCACATCGTAGCTACGATGCATTAGCAAGCTATCGGATCAGCTACCACATCGTAGC

Genom je razbijen na slučajno odabranim pozicijama

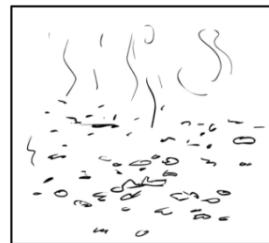


CTGATGATGGACTACGC ACTACTGCT GCTGTATTAGATCAGCTACC ATCGTAGCTAGATGCATTAGC AGCTATCG ATCAGCTAC CATCGTAGC
CTGATGATGGACTACGC CGCTACTACTCTAGCTGTATACGATCAGC ACCACATCGT GCTACGATGCA TAGCAAGCT CCGGATCAC TACCACAT GTAGC
CTGATGATGGACTACGC ACTACTGCTCTGTATTACATCAGCTACACATCGTAGC ACGATGCATTACAAGCTATGGATCAGCT CACATCGTAGC
CTGATGATGGACTACGC ACTGCTAGCTATTACGATGCTACCACACGTAGCTACGA GCATTAGCAA CTATCGGA CAGCTACCA ATCGTAGC

Generisana su očitavanja (*reads*)

CTGATGA TGGACTACGCTAC TACTGCTAG CTGTATTACG ATCAGCTACCACA TCGTAGCTACG ATGCATTAGCAA GCTATCGGA TCAGCTACCA CATCGTAGC
CTGATGATG GACTACGCT ACTACTGCTA GCTGTATTACG ATCAGCTACC ACATCGTAGCT ACGATGCATTA GCAAGCTATC GGATCAGCTAC CACATCGTAGC
CTGATGATGG ACTACGCTAC TACTGCTAGCT GTATTACGATC AGCTACCAC ATCGTAGCTACG ATGCATTAGCA AGCTATCGG A TCAGCTACCA CATCGTAGC
CTGATGATGGACT ACGCTACTACT GCTAGCTGTAT TACGATCAGC TACCACATCGT AGCTACGATGCA TTAGCAAGCT ATCGGATCA GCTACCACATC GTAGC

Neka očitavanja su nestala u eksploziji



CTGATGA TGGACTACGCTAC TACTGCTAG CTGTATTACG ATCAGCTACCACA TCGTAGCTACG ATGCATTAGCAA GCTATCGGA TCAGCTACCA CATCGTAGC
CTGATGATG GACTACGCT ACTACTGCTA GCTGTATTACG ATCAGCTACC ACATCGTAGCT ACGATGCATTA GCAAGCTATC GGATCAGCTAC CACATCGTAGC
CTGATGATGG ACTACGCTAC TACTGCTAGCT GTATTACGATC AGCTACCAC ATCGTAGCTACG ATGCATTAGCA AGCTATCGG A TCAGCTACCA CATCGTAGC
CTGATGATGGACT ACGCTACTACT GCTAGCTGTAT TACGATCAGC TACCACATCGT AGCTACGATGCA TTAGCAAGCT ATCGGATCA GCTACCACATC GTAGC

Ne znamo sa kojih pozicija su očitavanja došla

A collection of DNA sequence fragments of different lengths and orientations, all rotated diagonally upwards from left to right. The sequences include:

- ATCAGCTACCA
- TACTGCTAG
- CTGATGA
- ATGCATTAGCA
- CTGATGATG
- ACGCTACT
- ACATCGTAGCT
- TACTGCTAGCT
- ATCGATGGACT
- ATCAGCTACC
- TACTGCTAGCT
- GCAAGCTATC
- GAATACGCT
- ATCGGATCA
- GGATCAGCTAC
- ATCGTAGCTACG
- GCTAGCTACG
- ACTACTGCTA
- GCAAGCTATC
- ACTACGCTAC
- TAGTACGCTAC
- ATCAGCTACCA
- TAGCAAGCT
- GCTTACACATC
- ATCAGCTACCA
- TACGATGCTAC
- AGCTACAC
- GTATTACGATC
- AGCTATCGG
- TCGTAGCTAG
- CTGATGATGG
- ATGCATTAGCAA
- CACATCGTAGC
- TACACATCGT
- CTGATGATGG
- ATCGTAGCTACG
- ATCGTAC
- CATCGTAGC
- TCAGCTACCA
- CTGATGATGG
- AGCTACGATGCA
- AGCTATCGA

Ne znamo sa kojih pozicija su
očitavanja došla

Ne znamo sa kojih pozicija su očitavanja došla

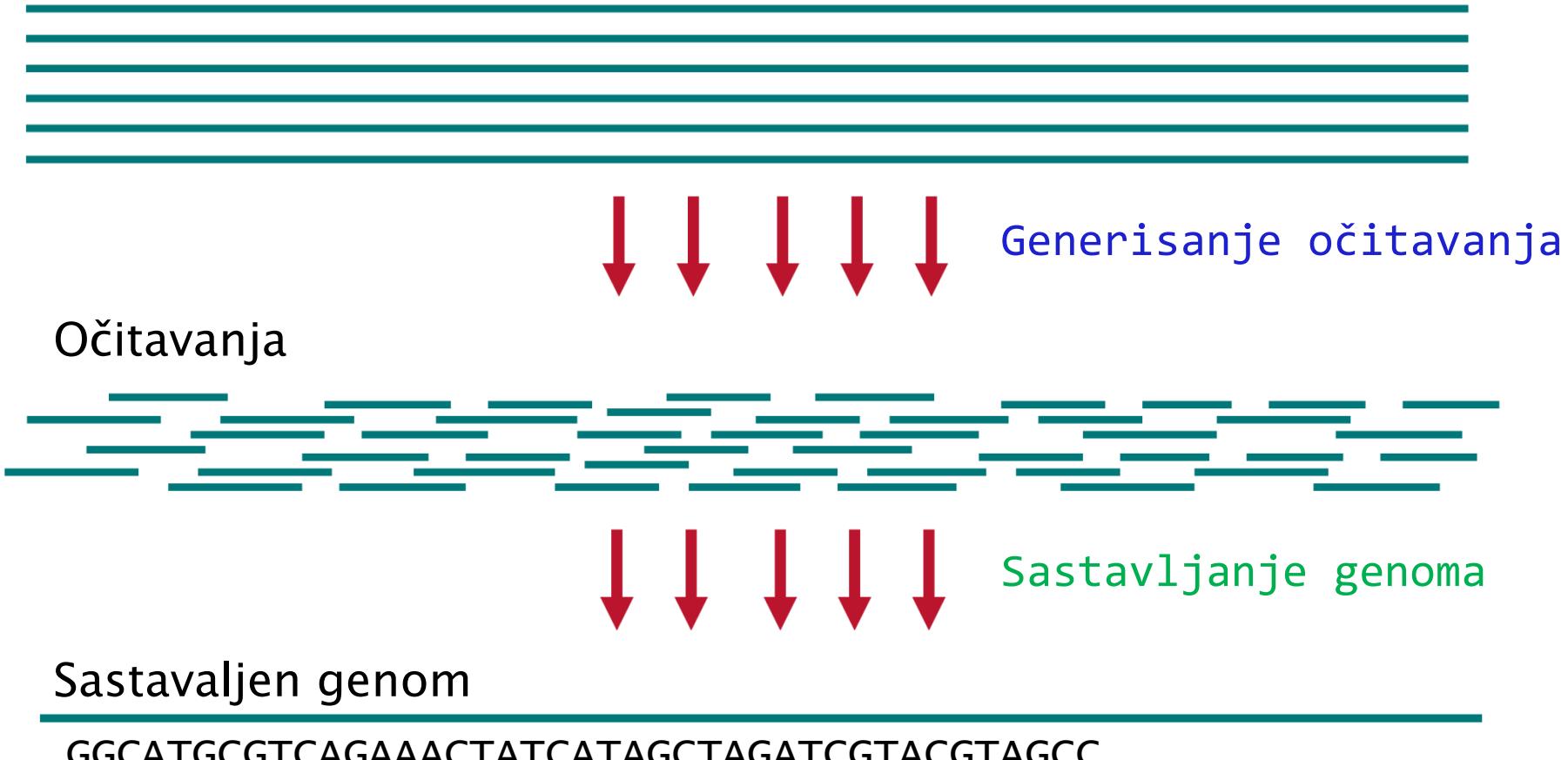
A collection of DNA sequence fragments arranged in a cloud-like pattern. The sequences are oriented at various angles. Two specific sequences are highlighted with yellow boxes: "GCTATCGGA" and "GCAAGCTATC".

The sequences include:

- ATCAGCTACCA
- TACTGCTAG
- CTGATGATGGACT
- ATCAGCTACC
- TGGACTACGCTAC
- AGCTATCGG
- AGCTACGATGCA
- ATGCATRGA
- CTGATGA
- TACTGCTAGCT
- GCTGTATTACG
- TTAGCAAGCT
- TCGTAGCTAG
- CTGATGATGG
- CATCGTAGC
- TCAGCTACCA
- ATGCATTAGCA
- TACTGCTAGCT
- GCTATCGGA
- GCAAGCTATC
- CTGTATTACG
- GCTACCACATC
- ATGCATTAGCAA
- ATGCATTAGCA
- CACATCGTAGC
- CTGATGATGG
- TACCATCGT
- ACGCTACT
- GGATCAGCTAC
- ACTACTGCTA
- TACGATCAGC
- ACGATGCATTA
- ACATCGTAGCT
- ATCGGATCA
- GAECTACGCT
- CTGTATTACG
- ATCGTAGCTACG
- AGCTACCAC
- CTGATGATGG
- ATCGTAGCTACG
- TACTGCTAGCT
- GCTAGCTGAT
- GTATTACGATC

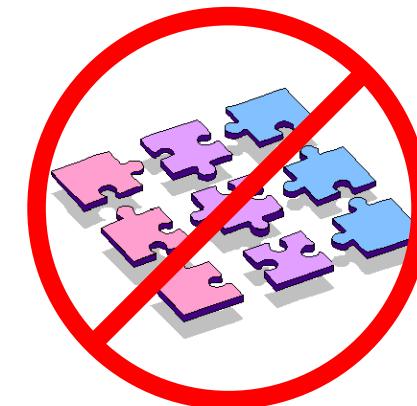
Od eksperimentalnih do računarskih problema

Više kopija genoma (nesekvencioniranog)



Zašto je sekvencioniranje genoma teško?

- Moderne mašine za sekvencioniranje (sekvenci) ne mogu da pročitaju ceo genom nukleotid po nukleotid od početka do kraja (kao što bismo pročitali knjigu)
- Mogu samo da iseckaju genom i generišu njegova kratka očitavanja
- Sastavljanje genoma nije isto kao i slaganje slagalice: moramo da koristimo preklapajuća očitavanja da bismo rekonstruisali genom



Pregled

- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- **Problem rekonstrukcije niske**
- Rekonstrukcija niske kao problem Hamiltonove putanje
- Rekonstrukcija niske kao problem Ojlerove putanje
- De Bruijinovi grafovi
- Ojlerova teorema
- Spajanje parova očitavanja
- U realnosti

Problem sekvencioniranja genoma

Problem sekvencioniranja genoma. Rekonstruisati genom na osnovu očitavanja.

- Ulaz. Kolekcija niski *Reads*.
- Izlaz. Niska *Genome* rekonstruisana na osnovu *Reads*.



Jovana Kovačević, Bioinformatika

k -gramski sastav niske

Composition₃(TAATGCCATGGGATGTT) =

TAA

AAT

ATG

TGC

GCC

CCA

CAT

ATG

TGG

GGG

GGA

GAT

ATG

TGT

GTT

k -gramski sastav niske

Composition₃(TAATGCCATGGGATGTT) =

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

=

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

leksikografski poredak

Rekonstrukcija niske na osnovu njenog k -gramskog sastava

Problem rekonstrukcije niske.

Rekonstruisati nisku na osnovu njenog k -gramskog sastava.

- Ulaz. Kolekcija k -grama.
- Izlaz. Niska *Genome* takva da je $\text{Composition}_k(\text{Genome})$ ekvivalentno kolekciji k -grama.

Naivni pristup

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

Naivni pristup

AAT

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

TAA

Naivni pristup

AAT

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

TAA

Naivni pristup

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

TAA
AAT

Naivni pristup

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

TAA
AAT

Naivni pristup

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

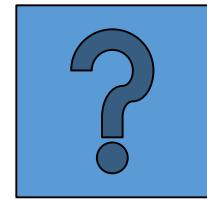
TAA
AAT
ATG

Naivni pristup

ATG ATG CAT CCA GAT GCC GGA GGG GTT

TGC TGG TGT

TA
AAT
ATG



Naivni pristup

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

TAA
AAT
ATG

Naivni pristup

ATG ATG CAT CCA GAT GCC GGA GGG **GTT** TGC TGG

TA**A**
A**AT**
ATG
TGT

Naivni pristup

ATG ATG CAT CCA GAT GCC GGA GGG **GTT** TGC TGG

TA**A**
A**AT**
ATG
TGT

Šta je sledeće?

ATG ATG CAT CCA GAT GCC GGA GGG

TGC TGG

TA^A
AAT
ATG
TGT
GTT



Pregled

- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- Problem rekonstrukcije niske
- **Rekonstrukcija niske kao problem Hamiltonove putanje**
- Rekonstrukcija niske kao problem Ojlerove putanje
- De Bruijinovi grafovi
- Ojlerova teorema
- Spajanje parova očitavanja
- U realnosti

Genom kao putanja

Composition₃(TAATGCCATGGGATGTT) =

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

Genom kao putanja

Composition₃(TAATGCCATGGGATGTT) =



Genom kao putanja

Composition₃(TAATGCCATGGGATGTT) =



Da li možemo konstruisati ovu **genomsku putanju** ako ne znamo sam genom TAATGCCATGGGATGTT ali znamo njegov *k*-gramske sastav?

Genom kao putanja

$\text{Composition}_3(\text{TAATGCCATGGGATGTT}) =$

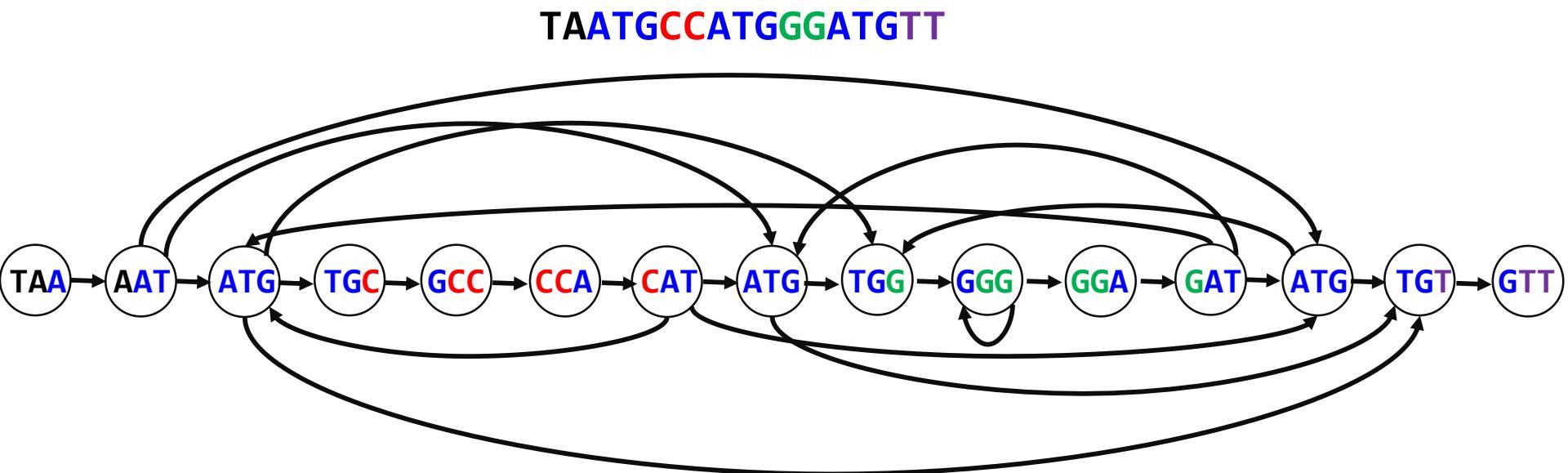


Da li možemo konstruisati ovu **genomsku putiju** ako ne znamo sam genom TAATGCCATGGGATGTT ali znamo njegov k -gramske sastav?

Možemo. Treba da povežemo $k\text{-mer}_1$ sa $k\text{-mer}_2$ ako $\text{suffix}(k\text{-mer}_1) = \text{prefix}(k\text{-mer}_2)$.

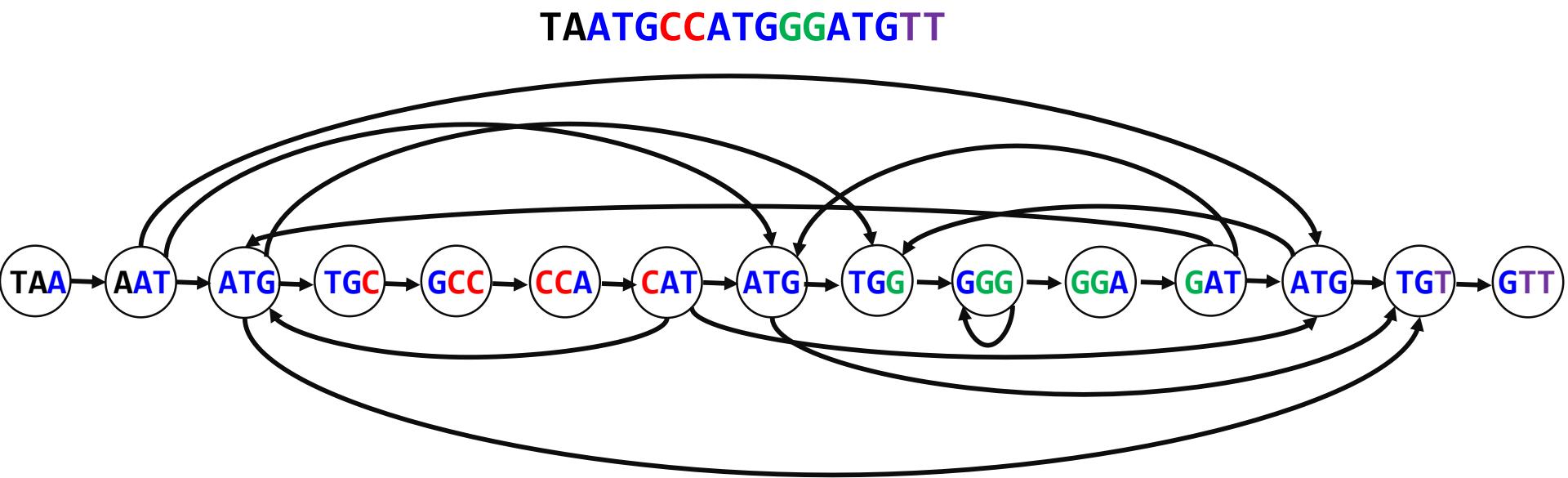
E.g. $\boxed{\text{TAA}} \rightarrow \boxed{\text{AAT}}$

Graf na osnovu k -gramskog sastava



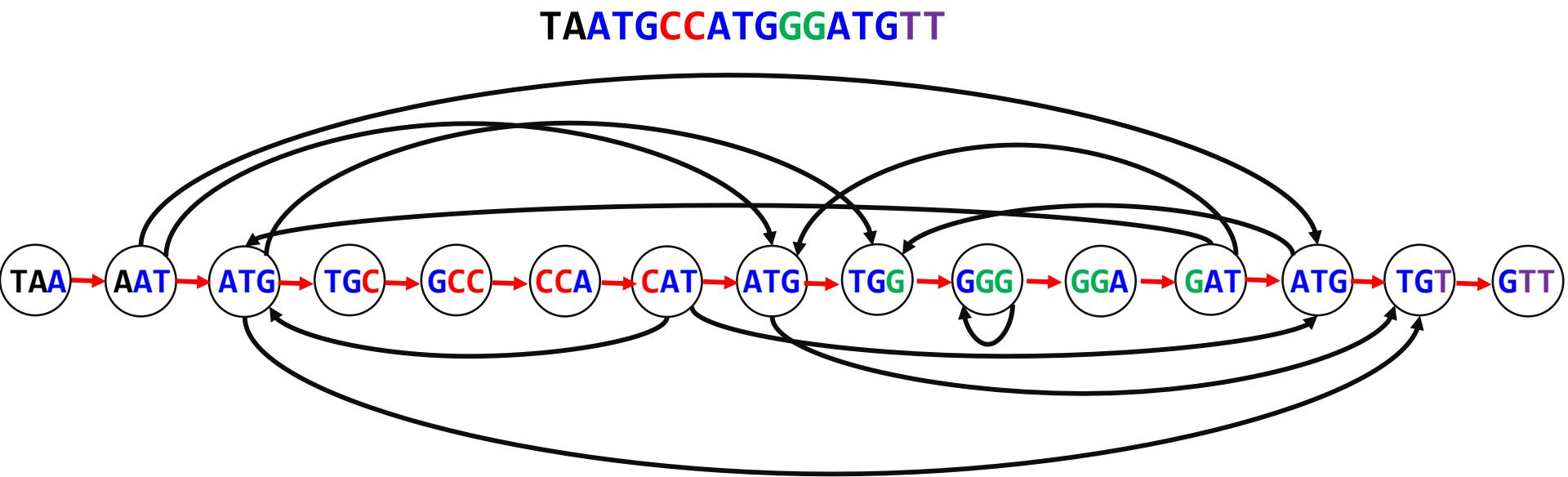
Možemo. Treba da povežemo k -mer₁ sa k -mer₂ ako
 $\text{suffix}(k\text{-mer}_1) = \text{prefix}(k\text{-mer}_2)$.
E.g. TAA → AAT

Graf na osnovu k -gramskog sastava



Od svih putanja, da li možemo da pronađemo
genomsку putanju u ovom grafu?

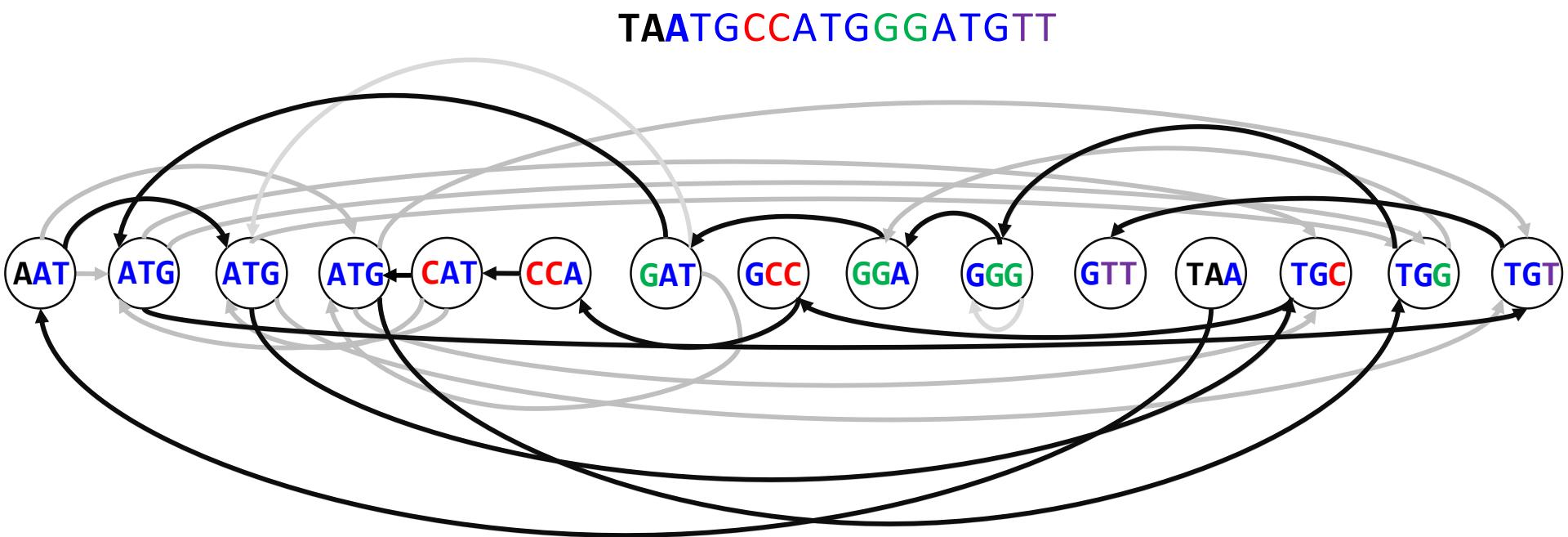
Graf na osnovu k -gramskog sastava



Od svih putanja, da li možemo da pronađemo genomsku putanju u ovom grafu?

Gde je genomska putanja?

Hamiltonova putanja: putanja koja posećuje svaki čvor u grafu tačno jednom.



Šta pokušavamo da pronađemo na ovom grafu?

Problem Hamiltonove putanje

Problem Hamiltonove putanje. Naći Hamiltonovu putanju u grafu.

- **Ulaz.** Graf.
- **Izlaz.** Putanja koja posećuje svaki čvor u grafu tačno jednom

Nalaženje Hamiltonove putanje je NP kompletan problem!

Pregled

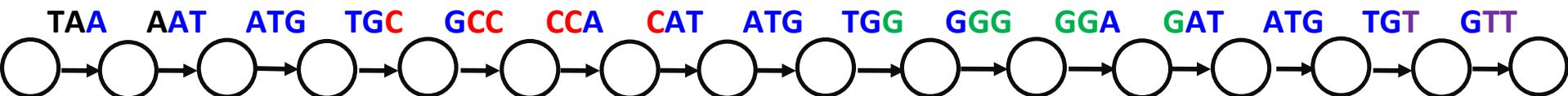
- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- Problem rekonstrukcije niske
- Rekonstrukcija niske kao problem Hamiltonove putanje
- **Rekonstrukcija niske kao problem Ojlerove putanje**
- De Bruijinovi grafovi
- Ojlerova teorema
- Spajanje parova očitavanja
- U realnosti

Malo drugačija putanja

TAATG**CC**ATGGGATGTT

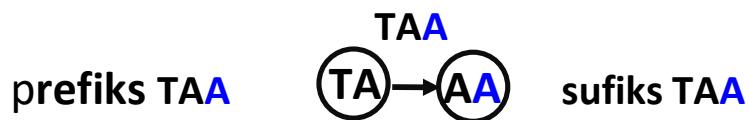


3-grami kao čvorovi



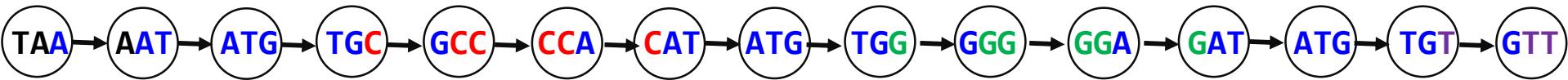
3-grami kao grane

Kako obeležavamo početni i krajnji čvor grane?

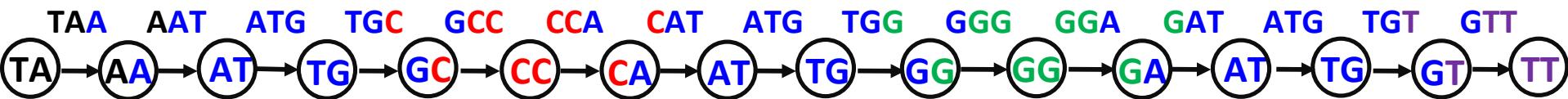


Obeležavanje čvorova u novoj putanji

TAATGCCATGGGATGTT

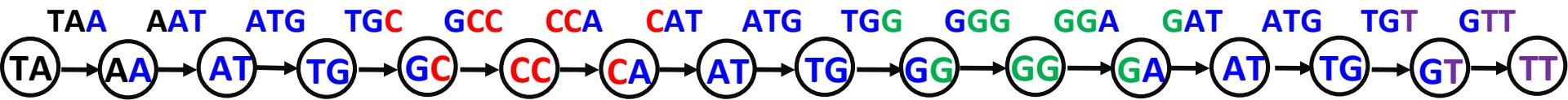


3-grami su čvorovi



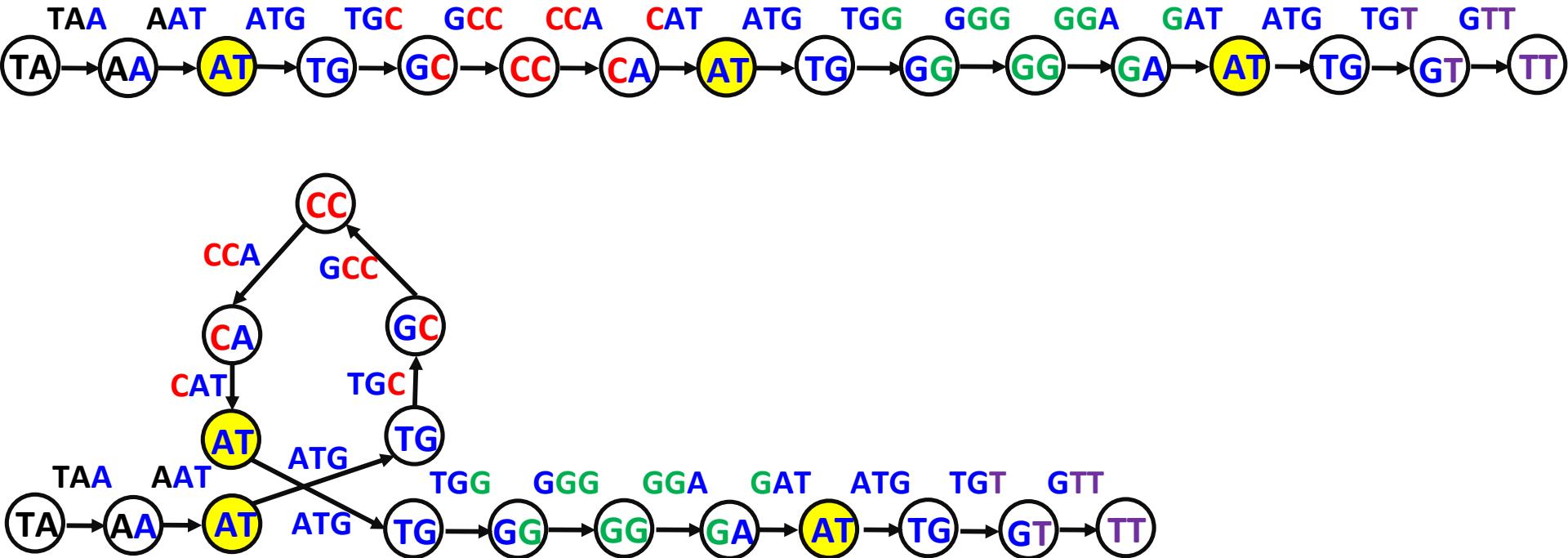
3-grami su grane a 2-grami su čvorovi

Obeležavanje čvorova u novoj putanji

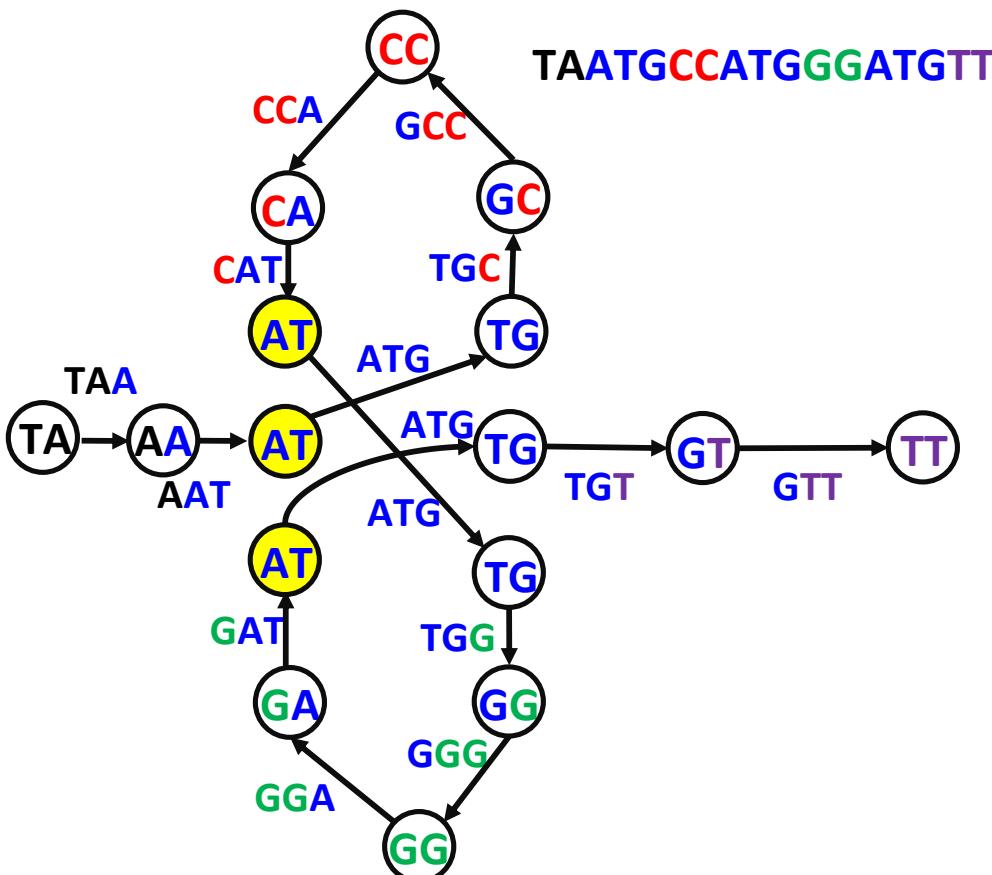


3-grami su grane a 2-grami su čvorovi

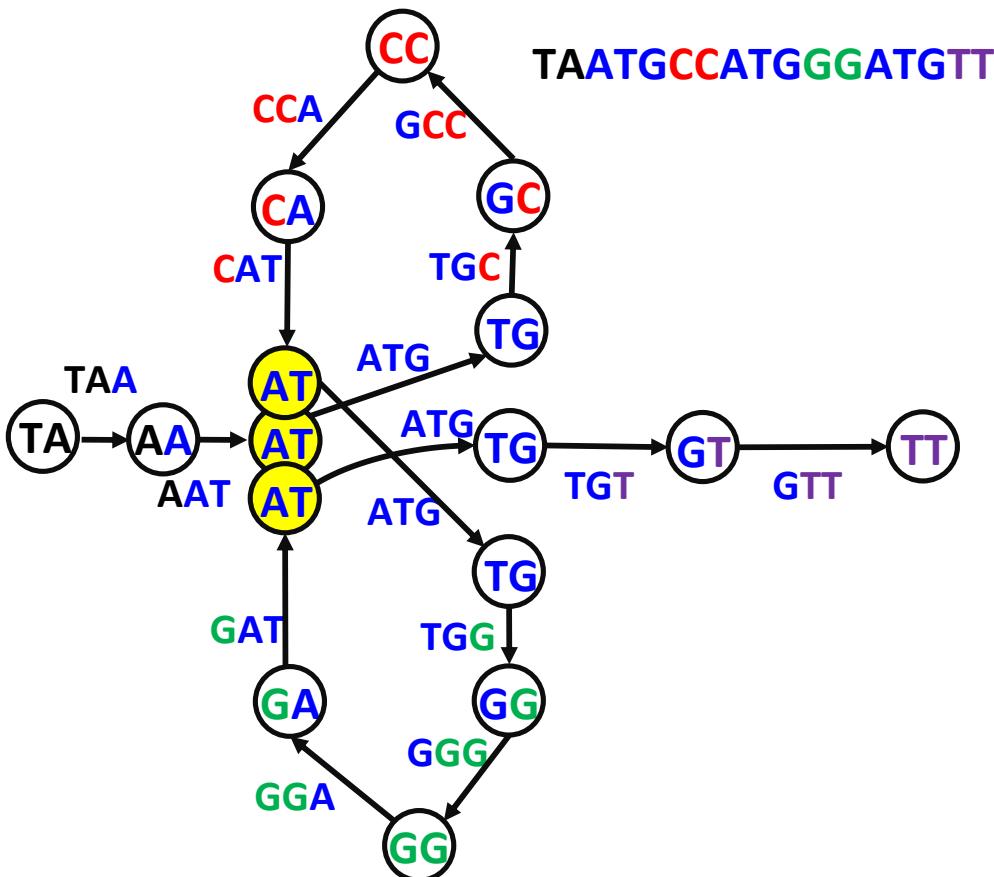
Lepljenje identično obeleženih čvorova



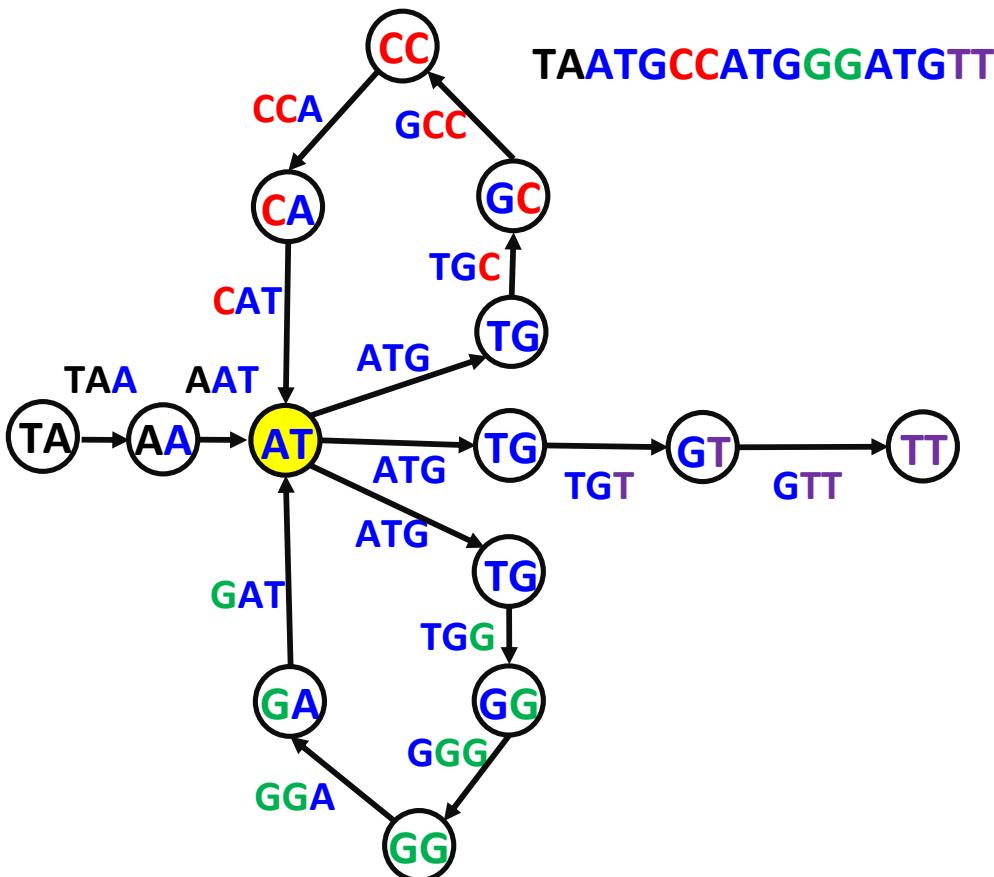
Lepljenje identično obeleženih čvorova



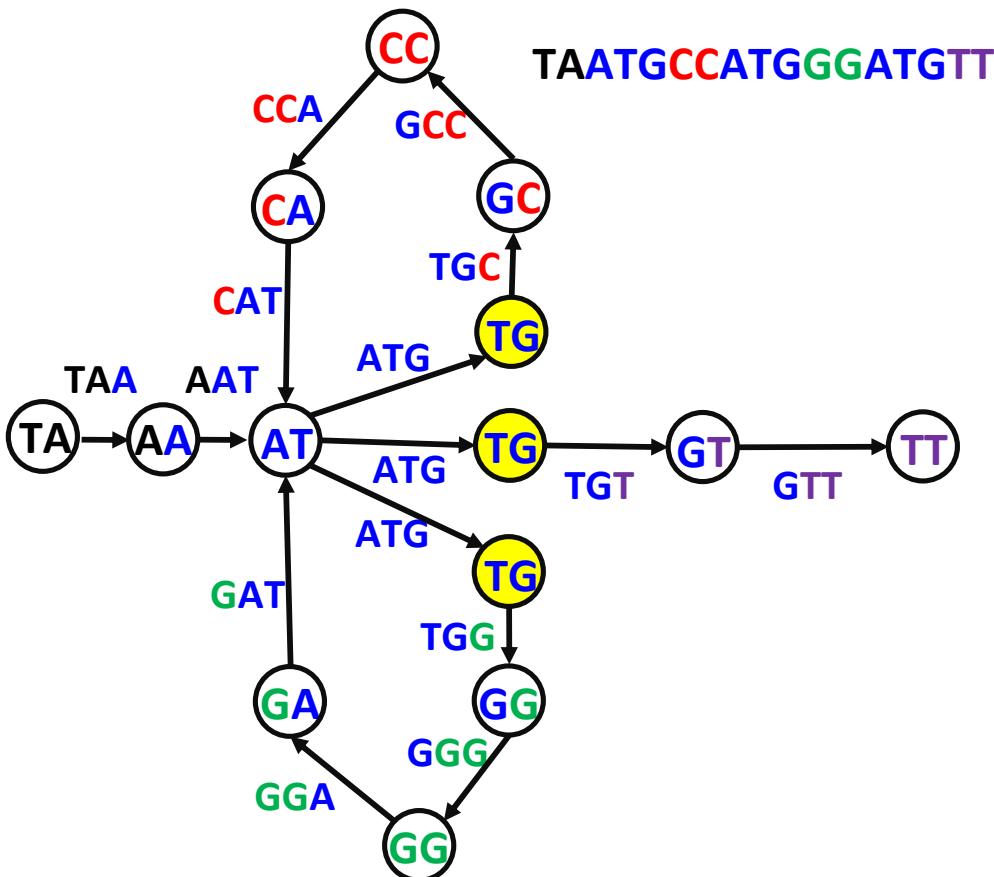
Lepljenje identično obeleženih čvorova



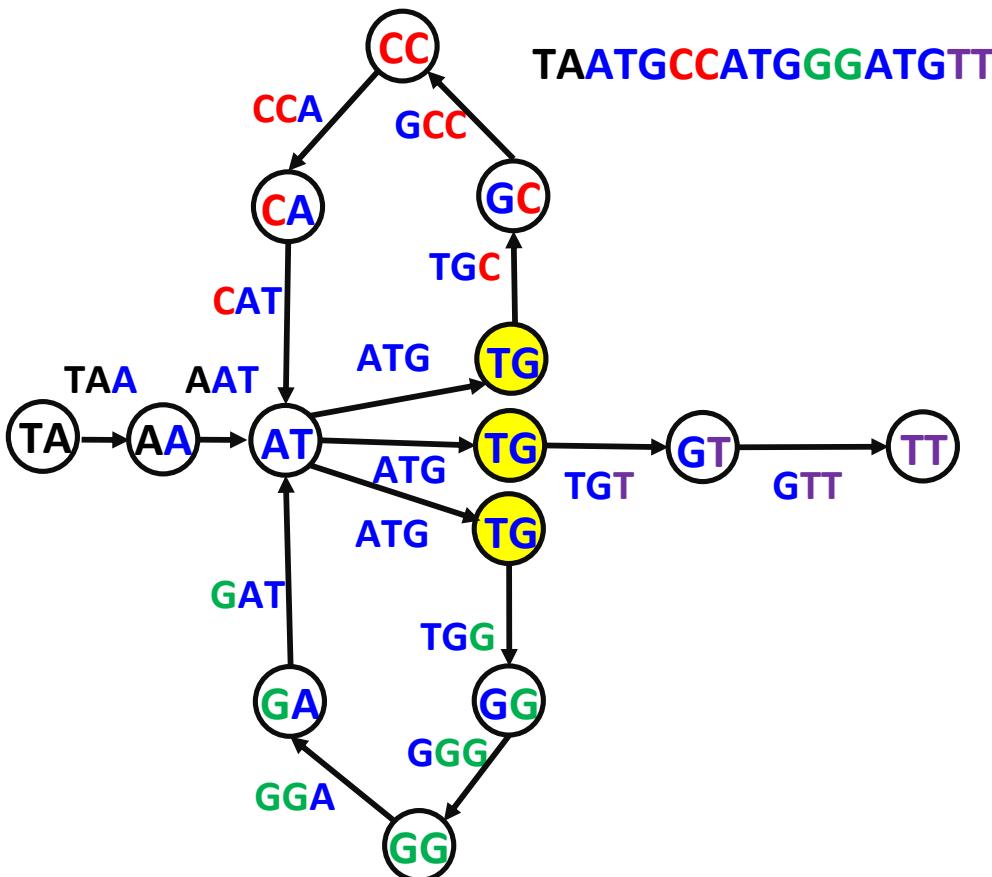
Lepljenje identično obeleženih čvorova



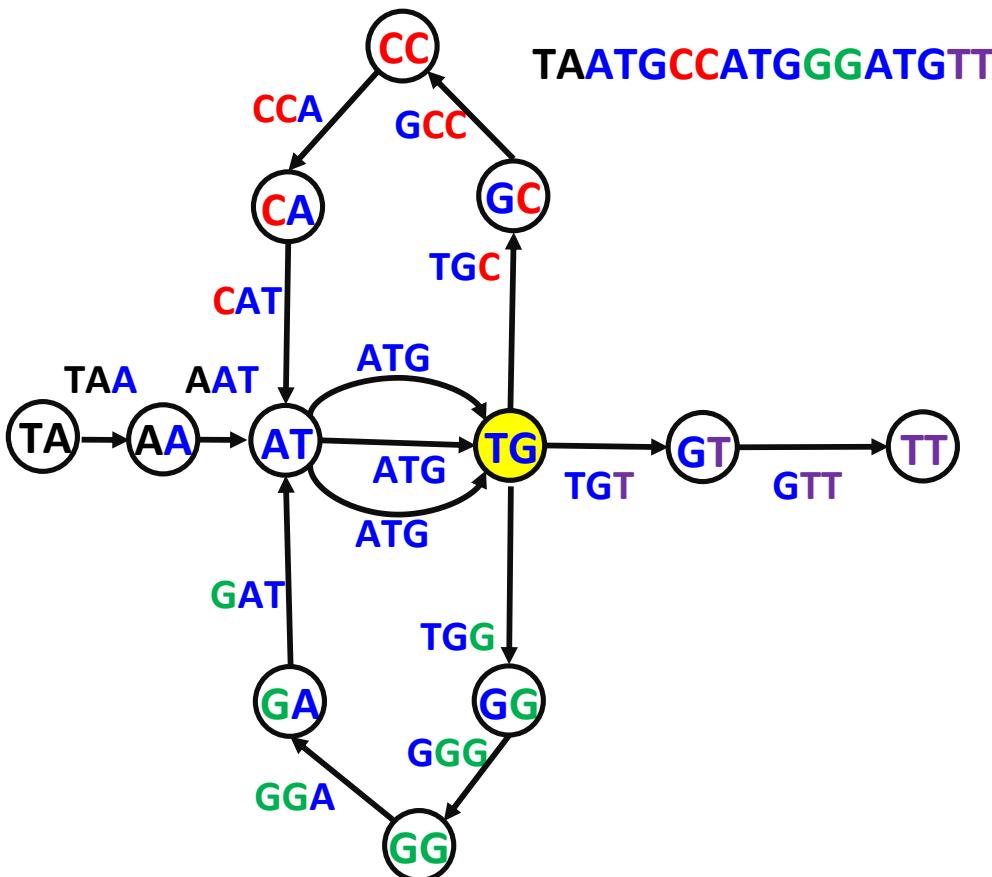
Lepljenje identično obeleženih čvorova



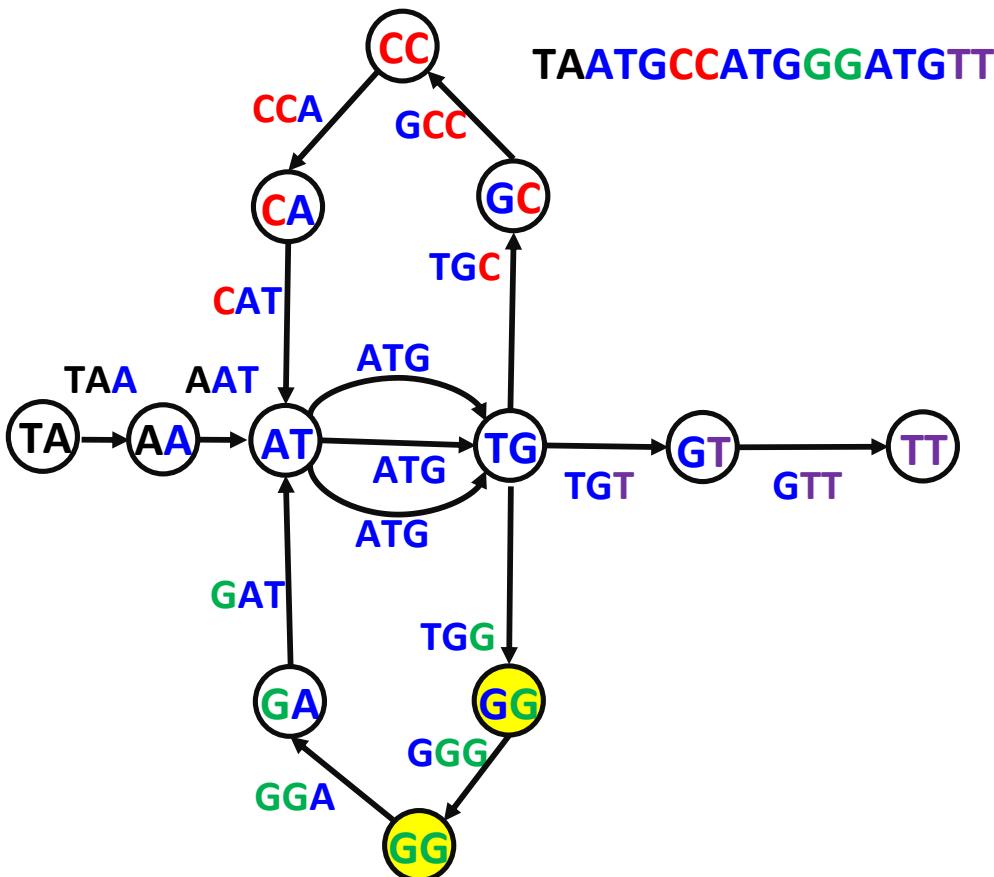
Lepljenje identično obeleženih čvorova



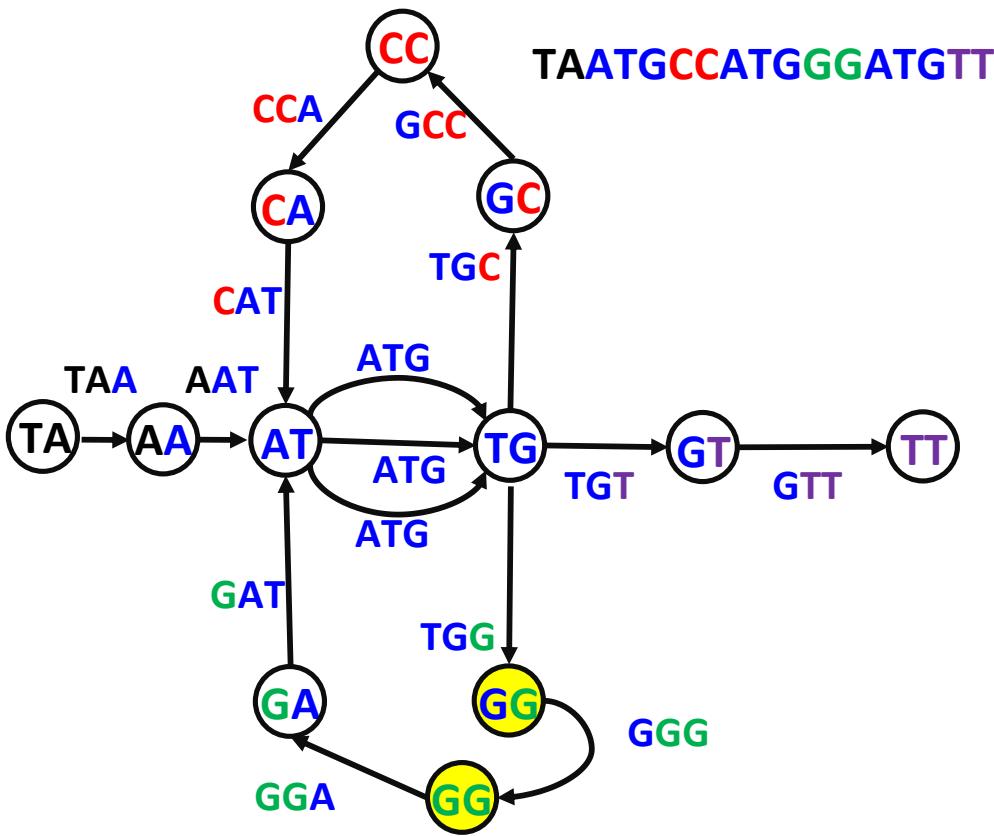
Lepljenje identično obeleženih čvorova



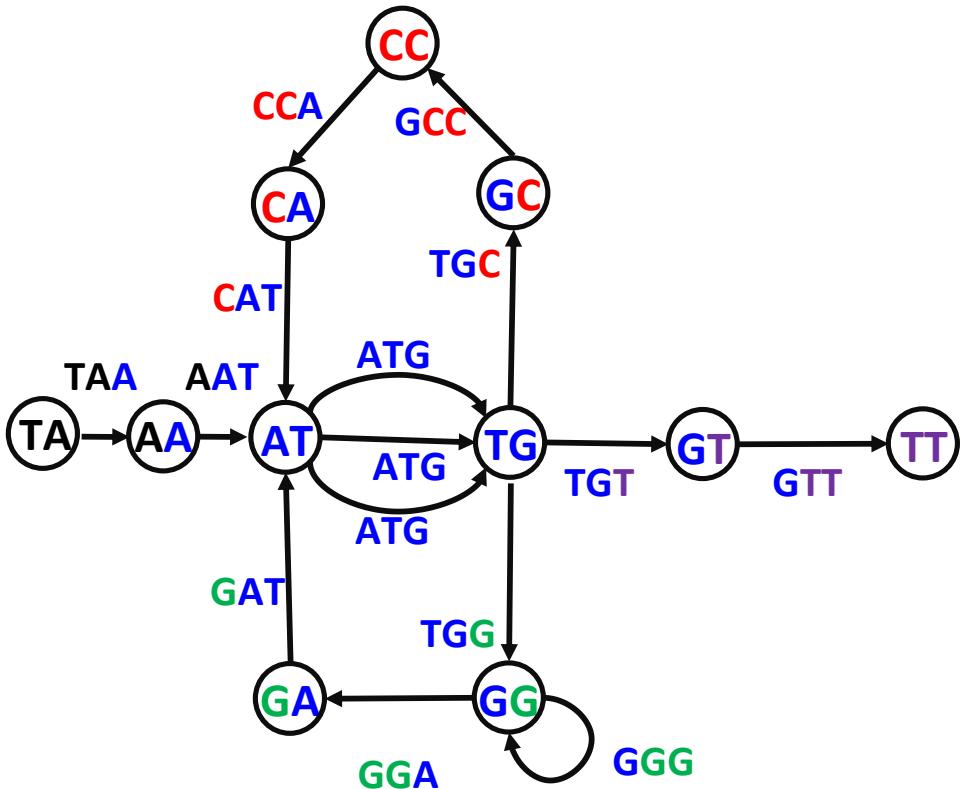
Lepljenje identično obeleženih čvorova



Lepljenje identično obeleženih čvorova



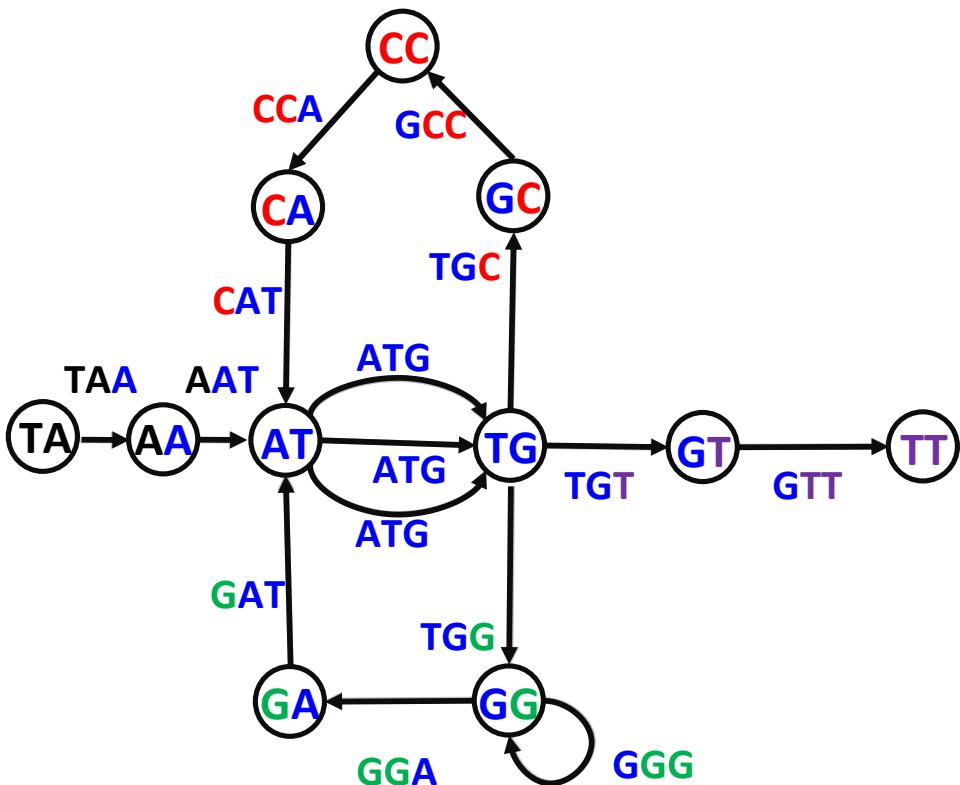
De Bruijinov graf za nisku TAATGCCATGGGATGTT



Gde se *Genome*
krije u ovom
grafu?

Gde je *Genome* u De Brujinovom grafu?

TAATGCCATGGGATGTT

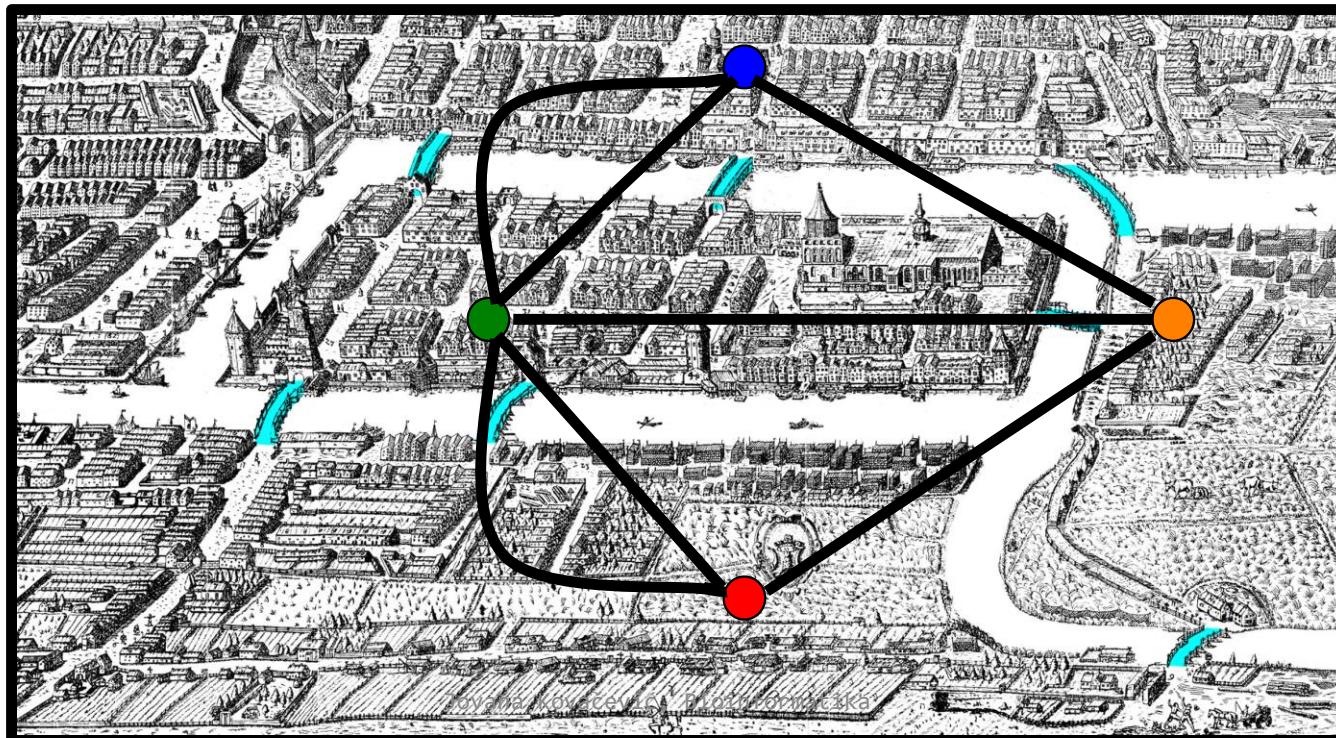


Ojlerova putanja u grafu je putanja koja posećuje svaku granu tačno jednom.

Problem Ojlerove putanje

Problem Ojlerove putanje. Pronaći Ojlerovu putanju u grafu.

- **Ulaz.** Graf.
- **Izlaz.** Putanja koja posećuje svaku granu u grafu tačno jednom.

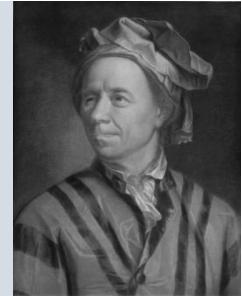


Pregled

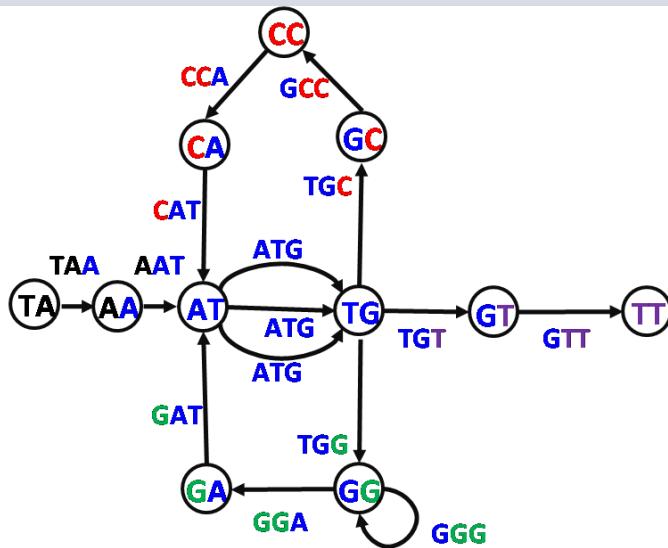
- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- Problem rekonstrukcije niske
- Rekonstrukcija niske kao problem Hamiltonove putanje
- Rekonstrukcija niske kao problem Ojlerove putanje
- **De Bruijinovi grafovi**
- Ojlerova teorema
- Spajanje parova očitavanja
- U realnosti

Problem Ojlerove putanje

Problem Ojlerove putanje. Pronaći Ojlerovu putanju u grafu.

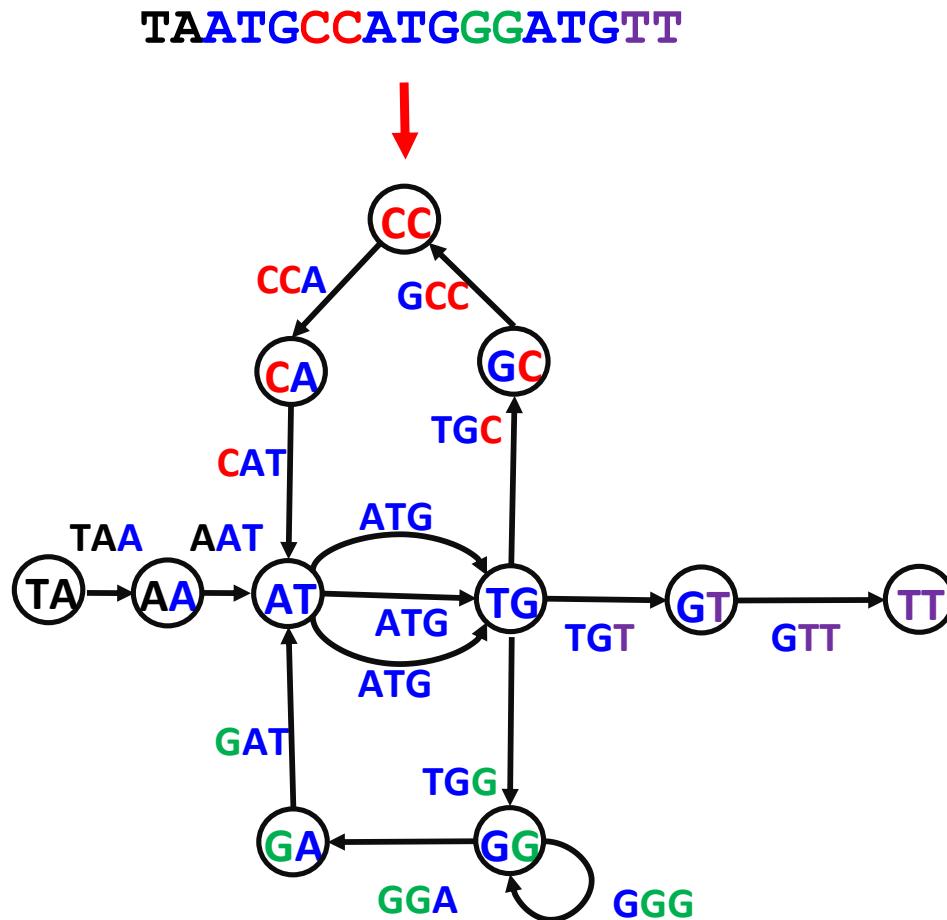


- **Ulaz.** Graf.
- **Izlaz.** Putanja koja posećuje svaku granu u grafu tačno jednom.



Konstruisali smo De Brujinov graf na osnovu genoma, ali u realnim primenama, genom je nepoznat!

Urađeno: Od genoma do De Bruijinovog grafa



Želimo da uradimo: Od očitavanja (kolekcije k -grama) do genoma

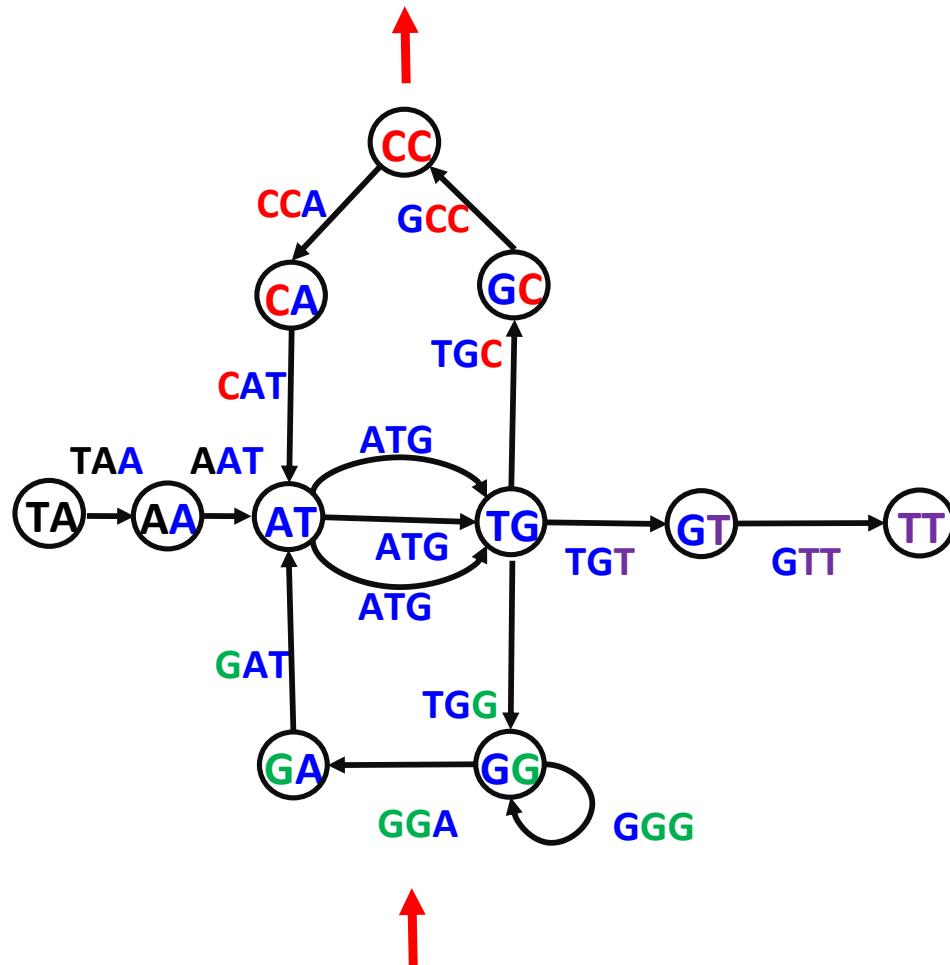
TAATGCCATGGGATGTT



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

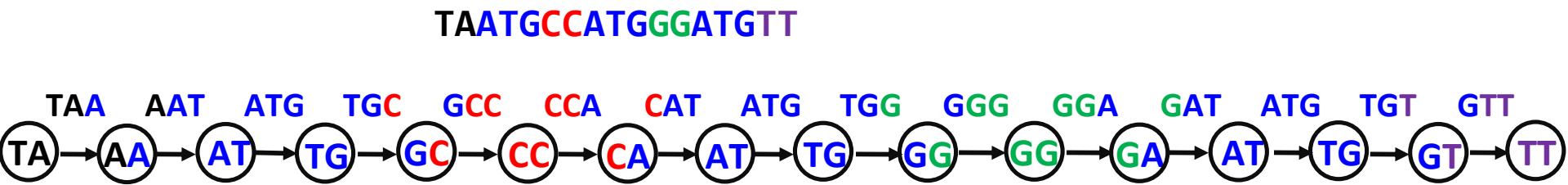
Pokazaćemo: Od očitavanja do De Brujinovog grafa do genoma

TA**ATGCCATGGGATGTT**



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

Konstrukcija De Brujinovog grafa kada je genom poznat



Konstrukcija De Brujinovog grafa kada je genom nepoznat

TAA ATG GCC CAT TGG GGA ATG GTT

AAT TGC CCA ATG GGG GAT TGT

Composition₃(TAATGCCATGGGATGTT)

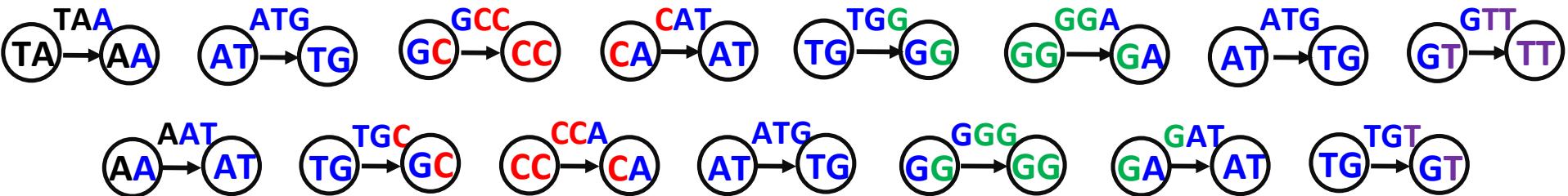
Predstavimo k -gramske sastave kao graf koji se sastoje od nepovezanih grana

TAA → ATG → GCC → CAT → TGG → GGA → ATG → GTT →

AAT → TGC → CCA → ATG → GGG → GAT → TGT →

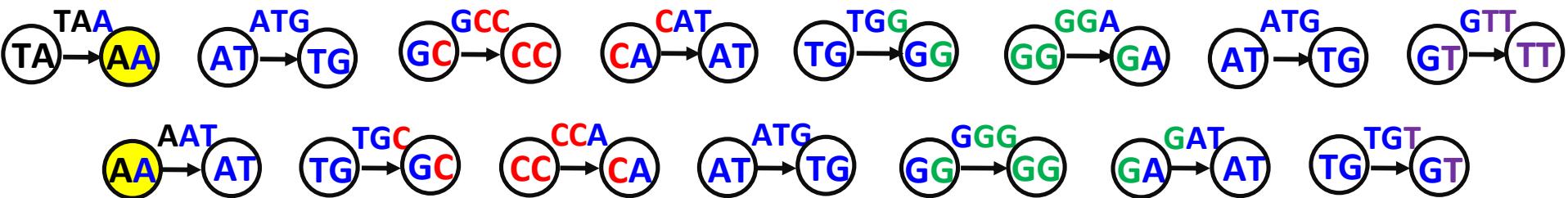
Composition₃(TAATGCCATGGGATGTT)

Konstruišemo De Brujinov graf na osnovu k -gramskog sastava

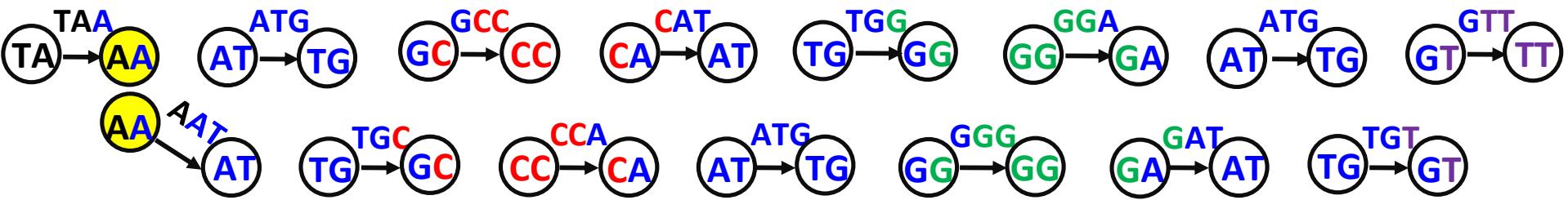


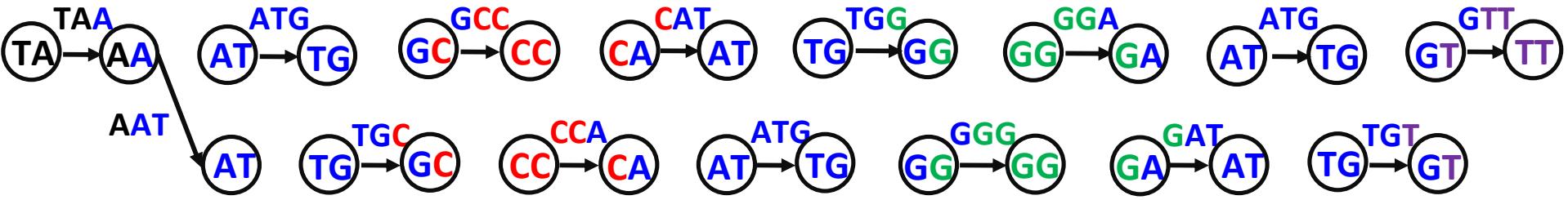
Composition₃(TAATGCCATGGGATGTT)

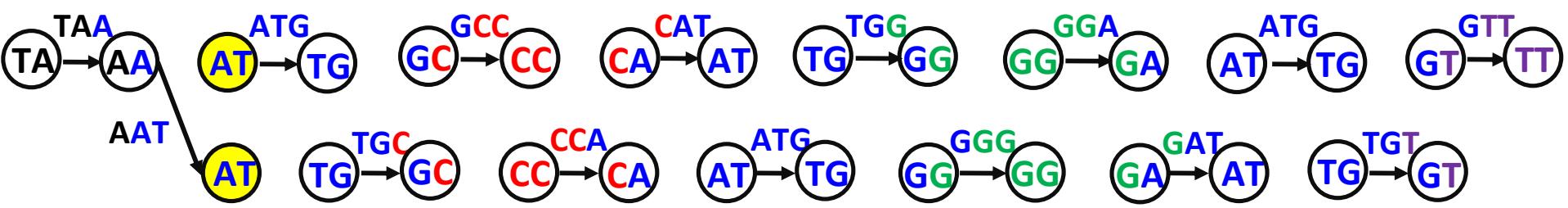
Zalepimo identično obeležene čvorove

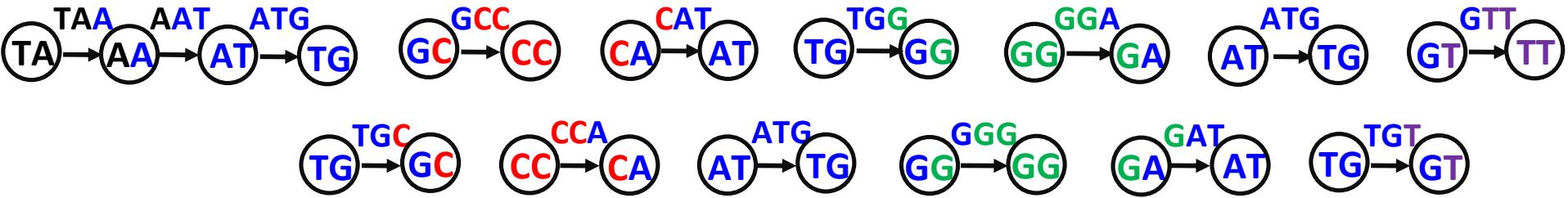


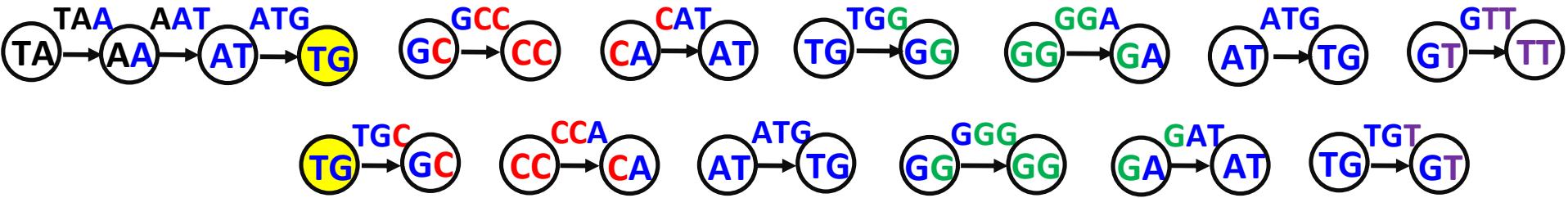
Zalepimo identično obeležene čvorove

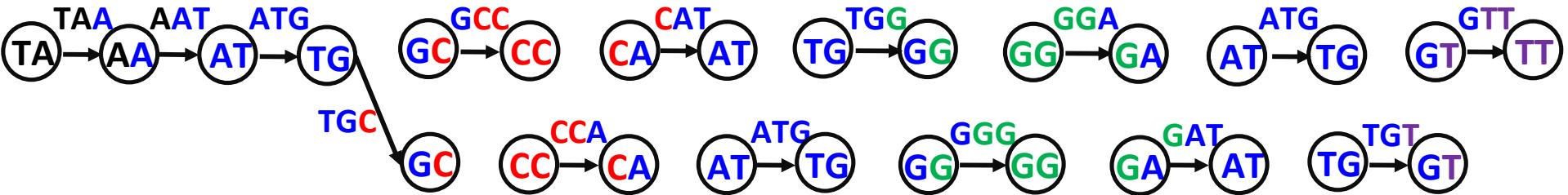


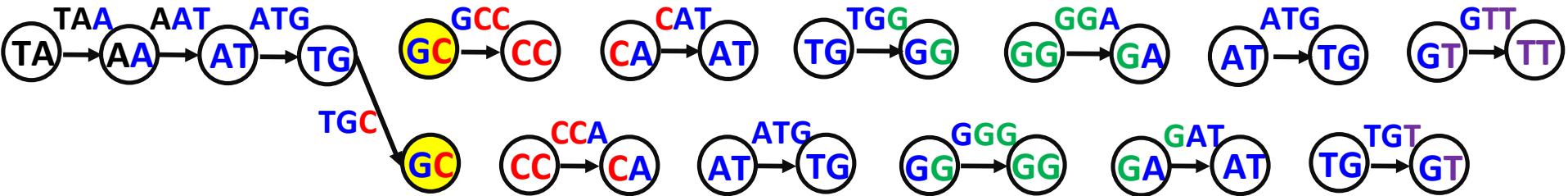


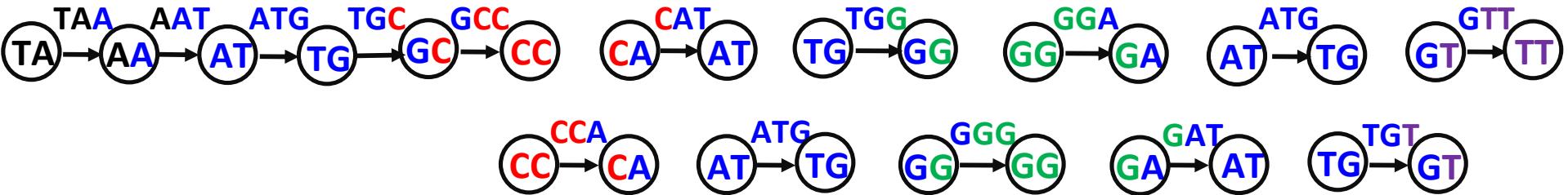


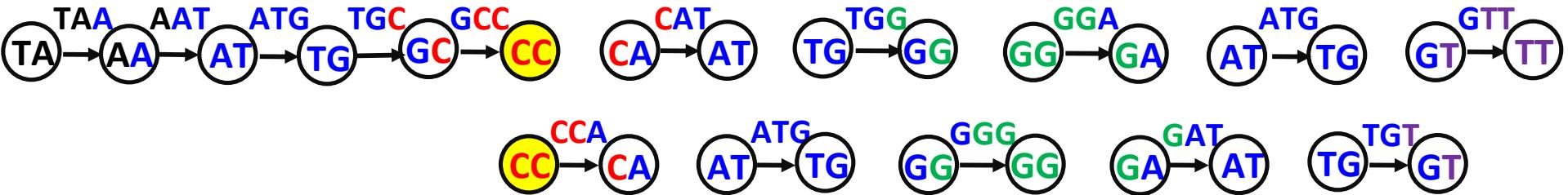


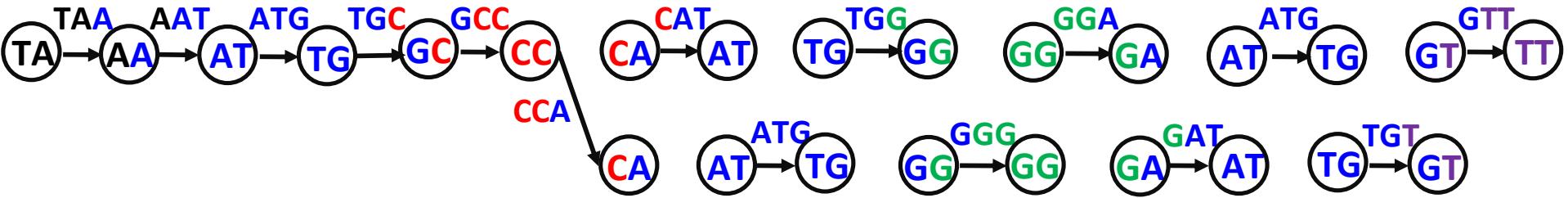


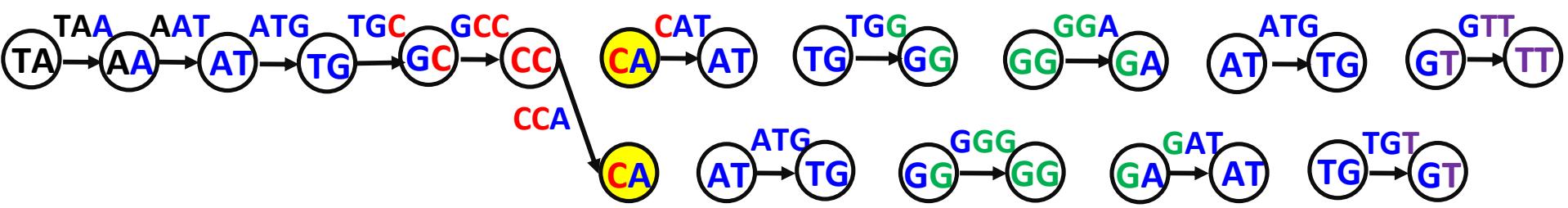


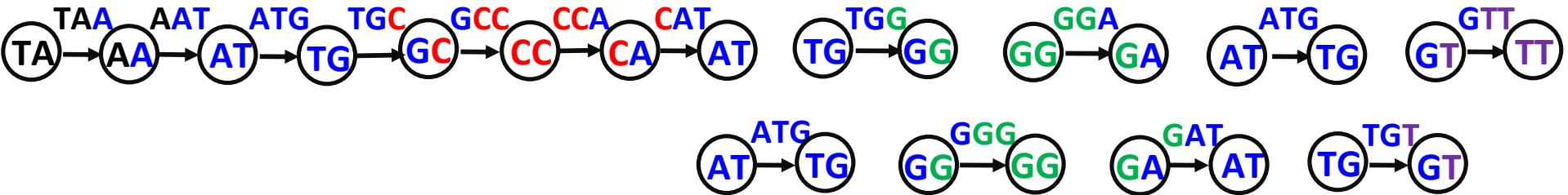


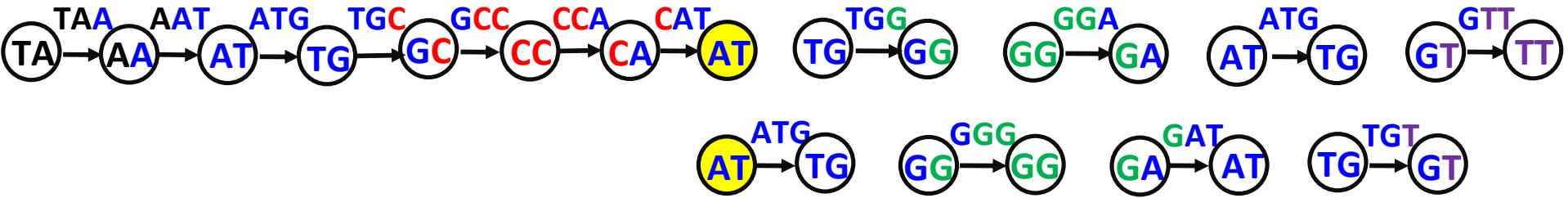


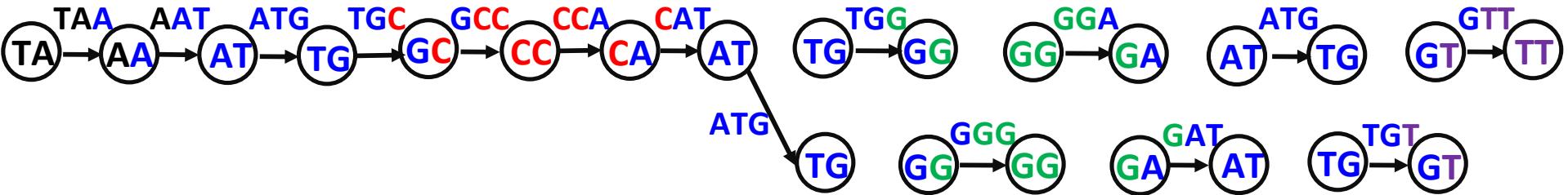


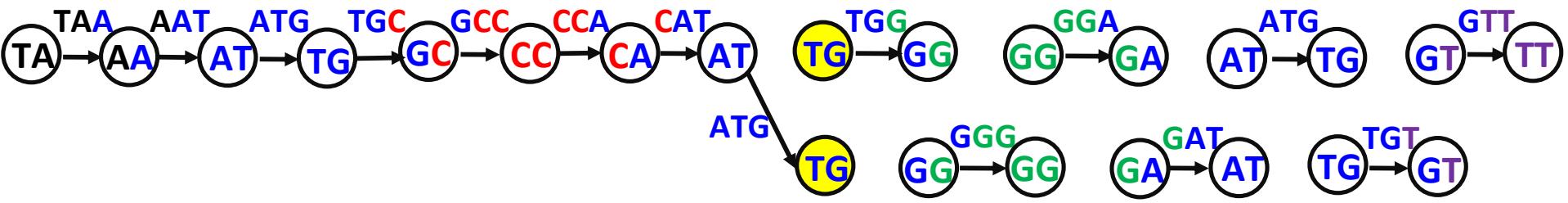


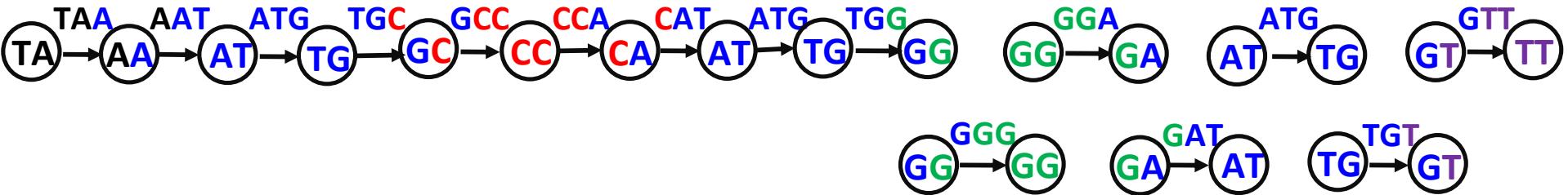


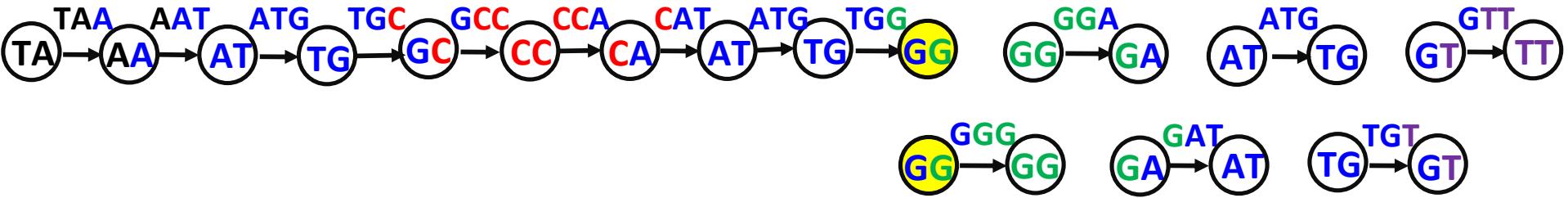


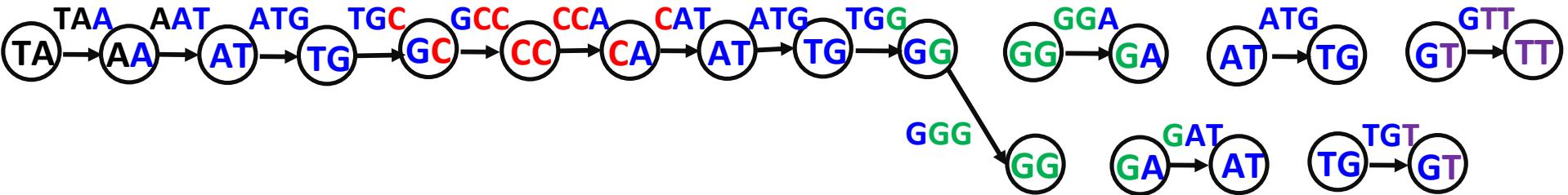


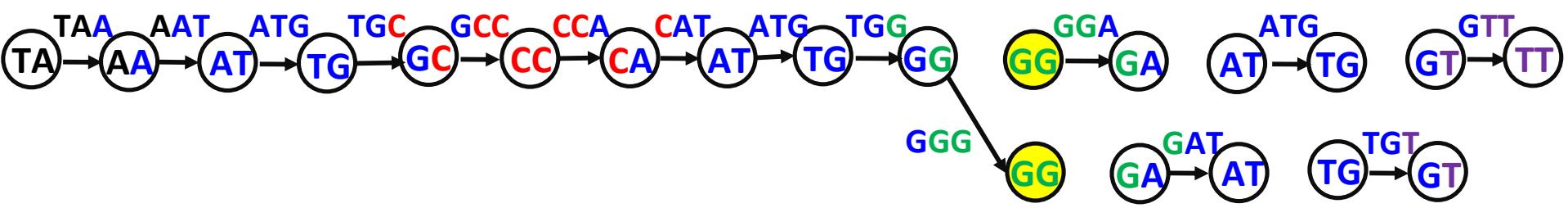


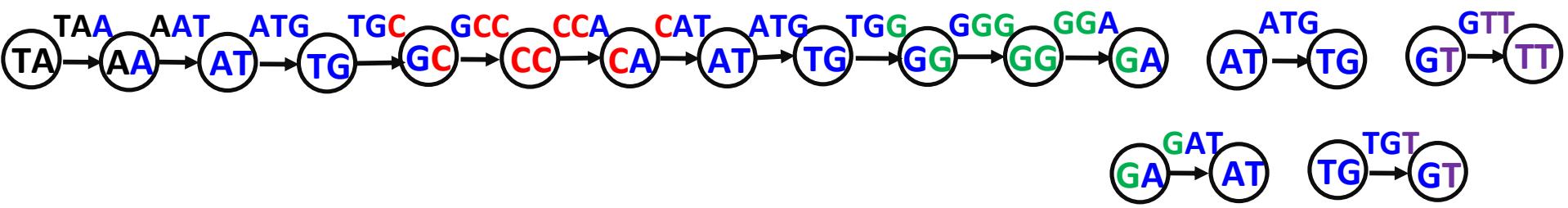


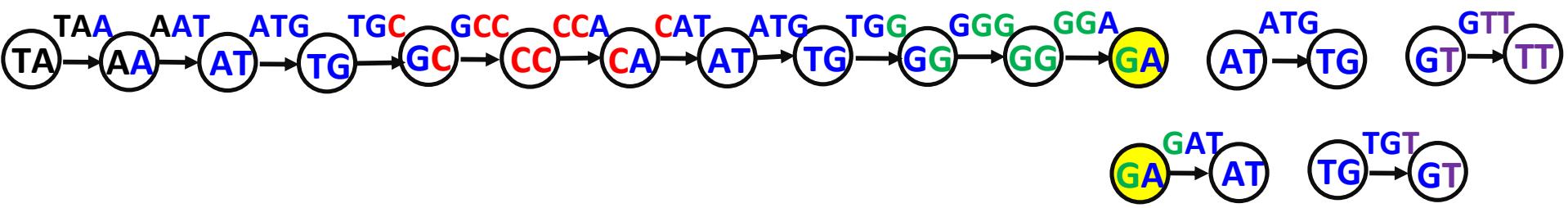


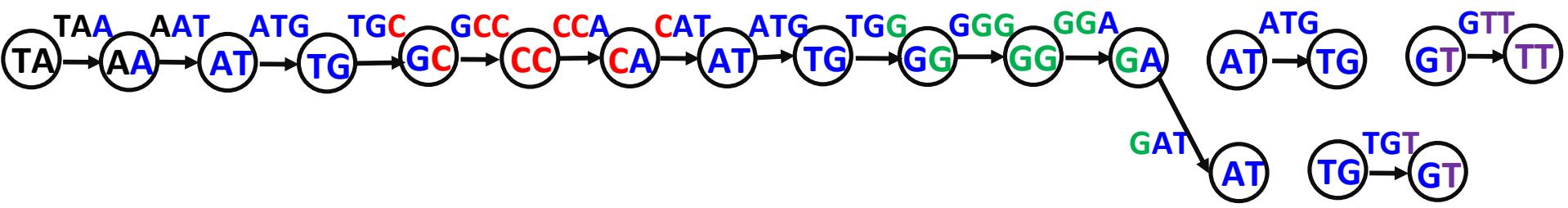


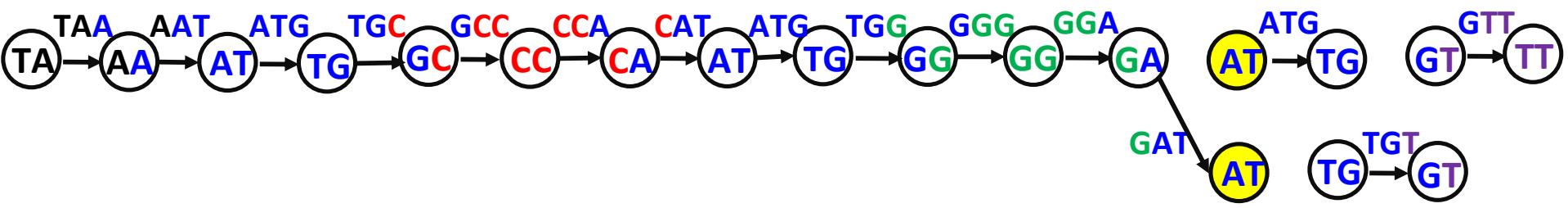


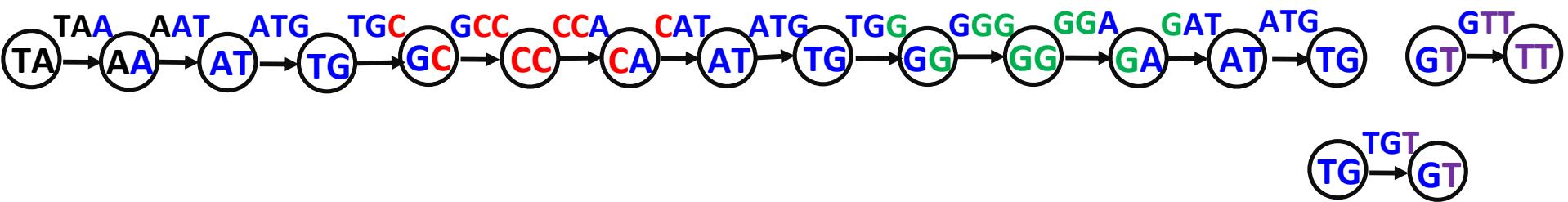


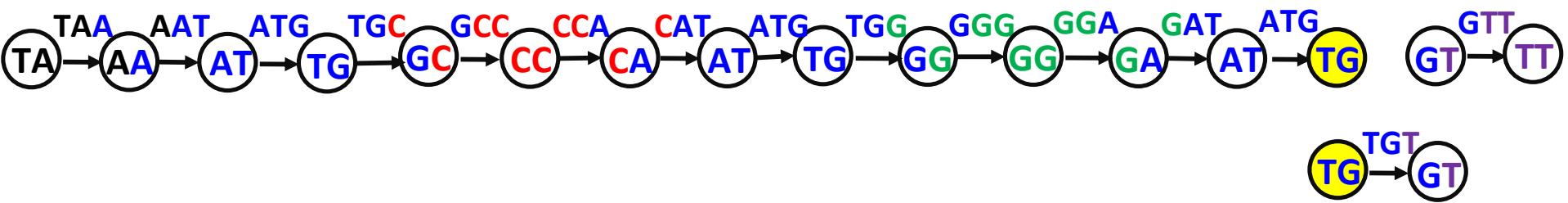


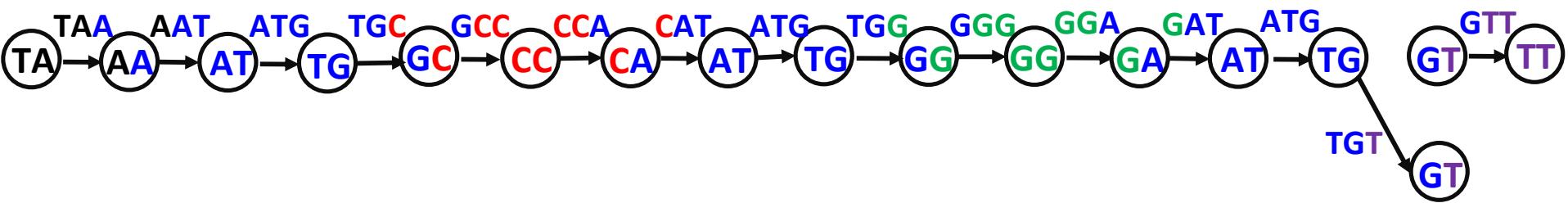


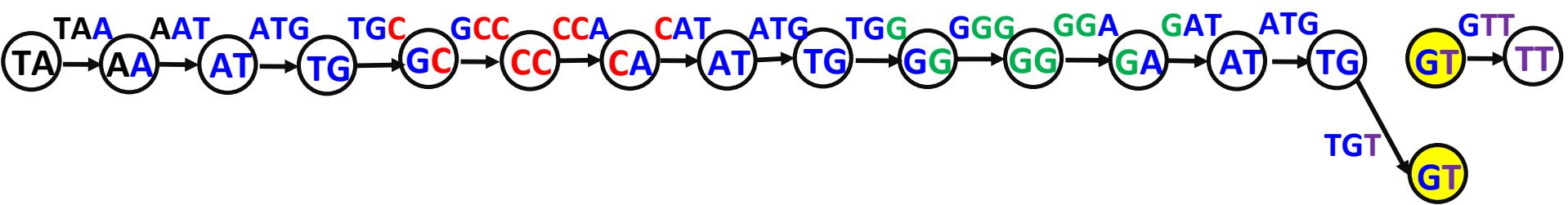




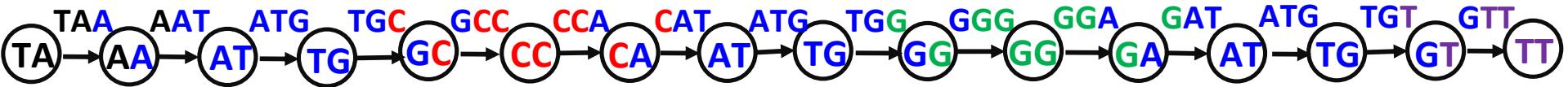




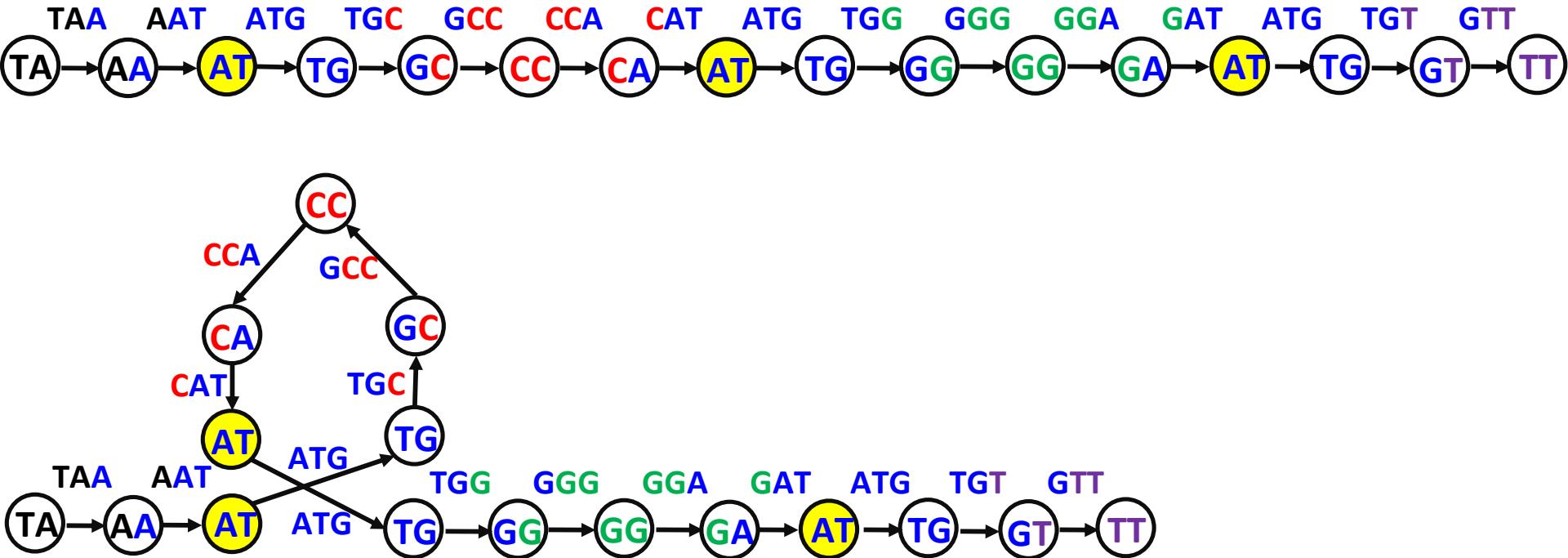




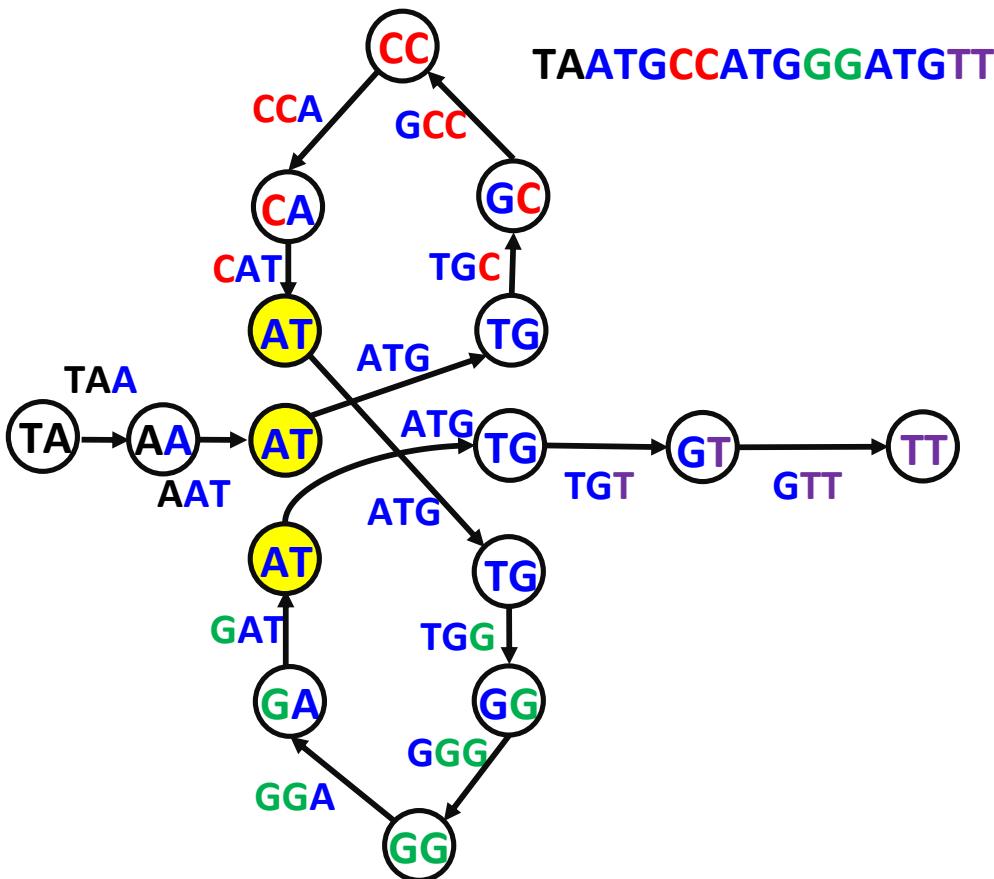
Lepljenje nije završeno

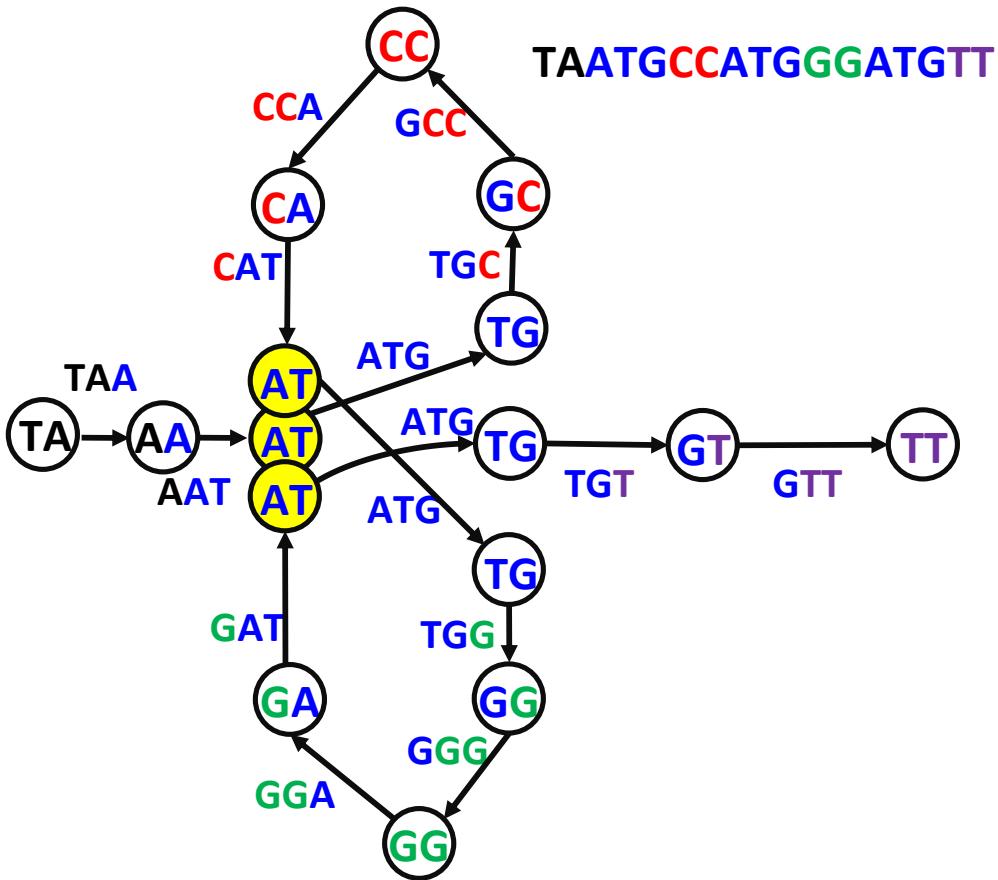


Zalepimo identično obeležene čvorove

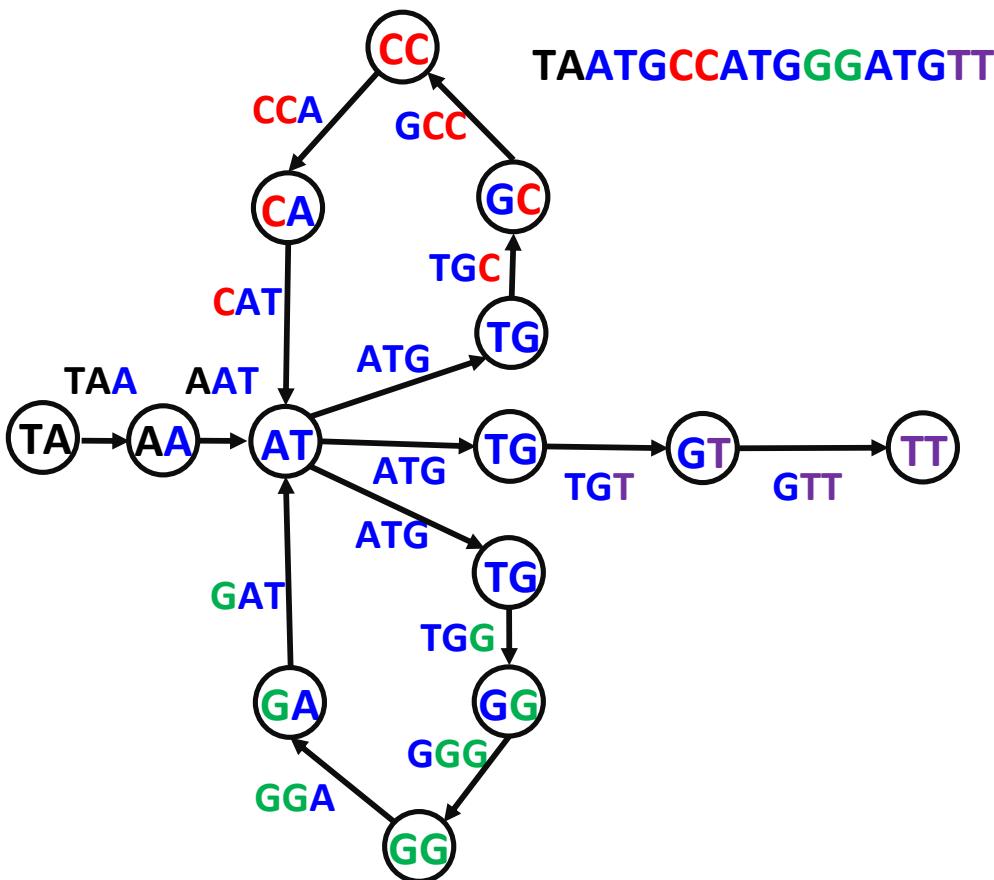


Zalepimo identično obeležene čvorove

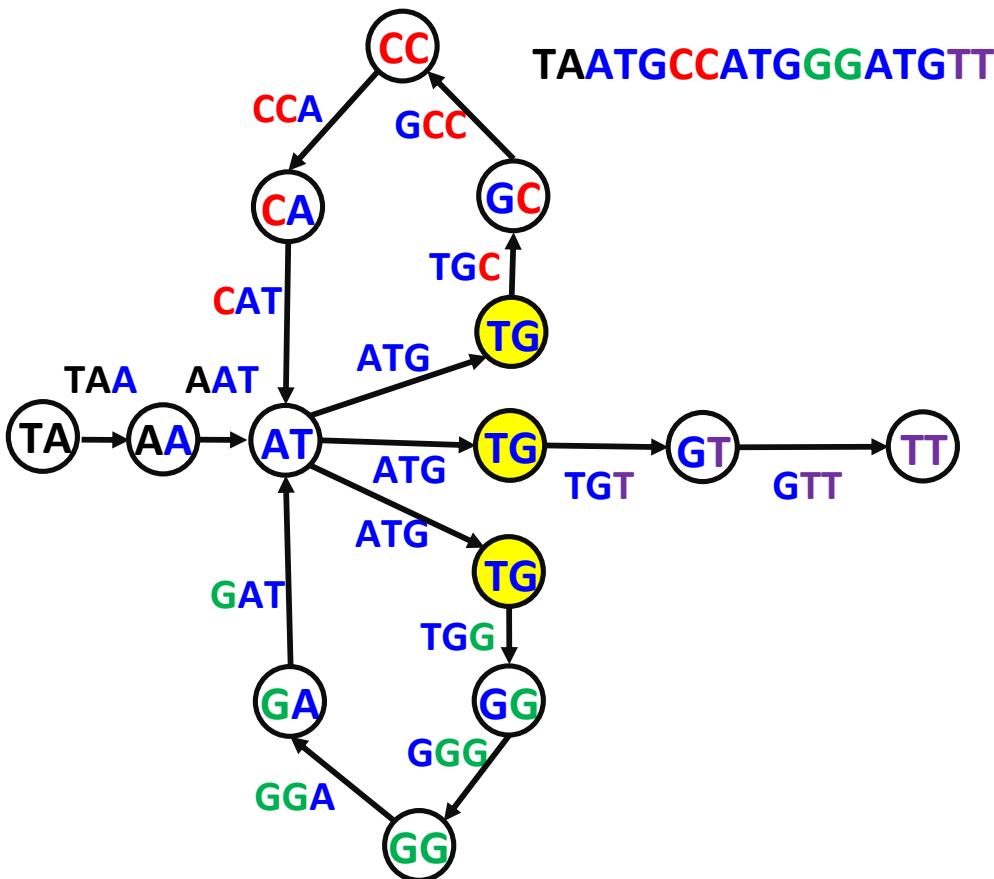




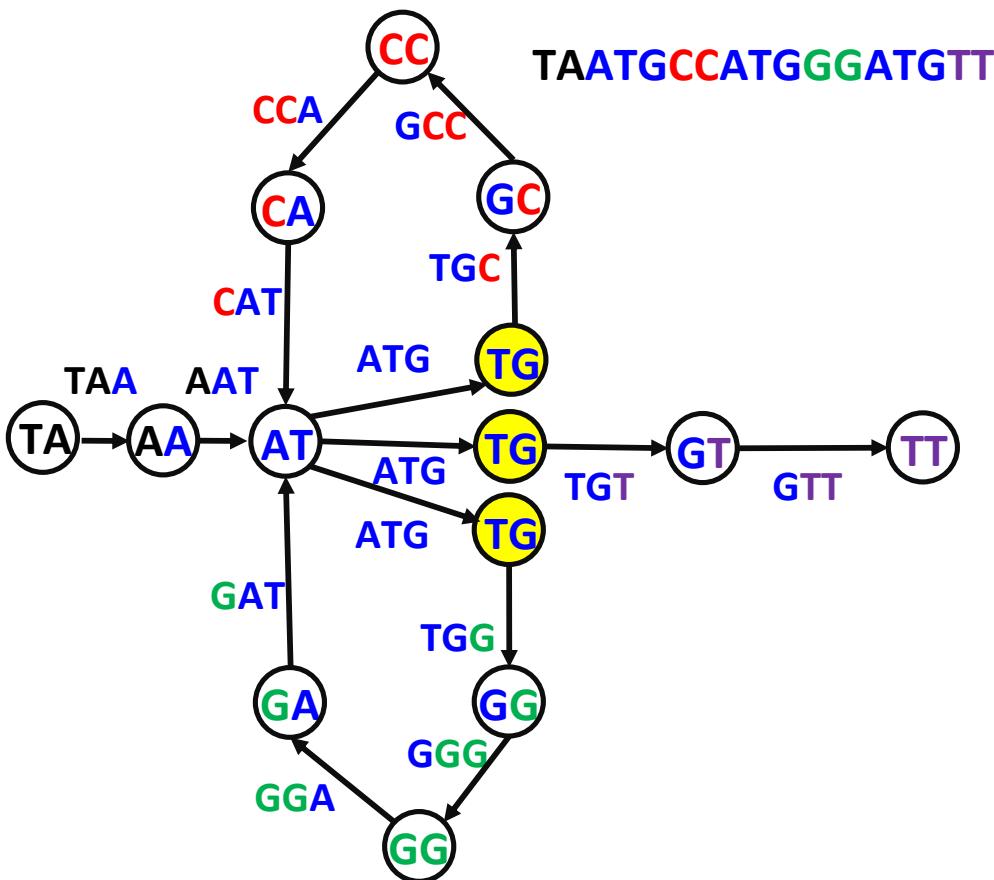
Zalepimo identično obeležene čvorove



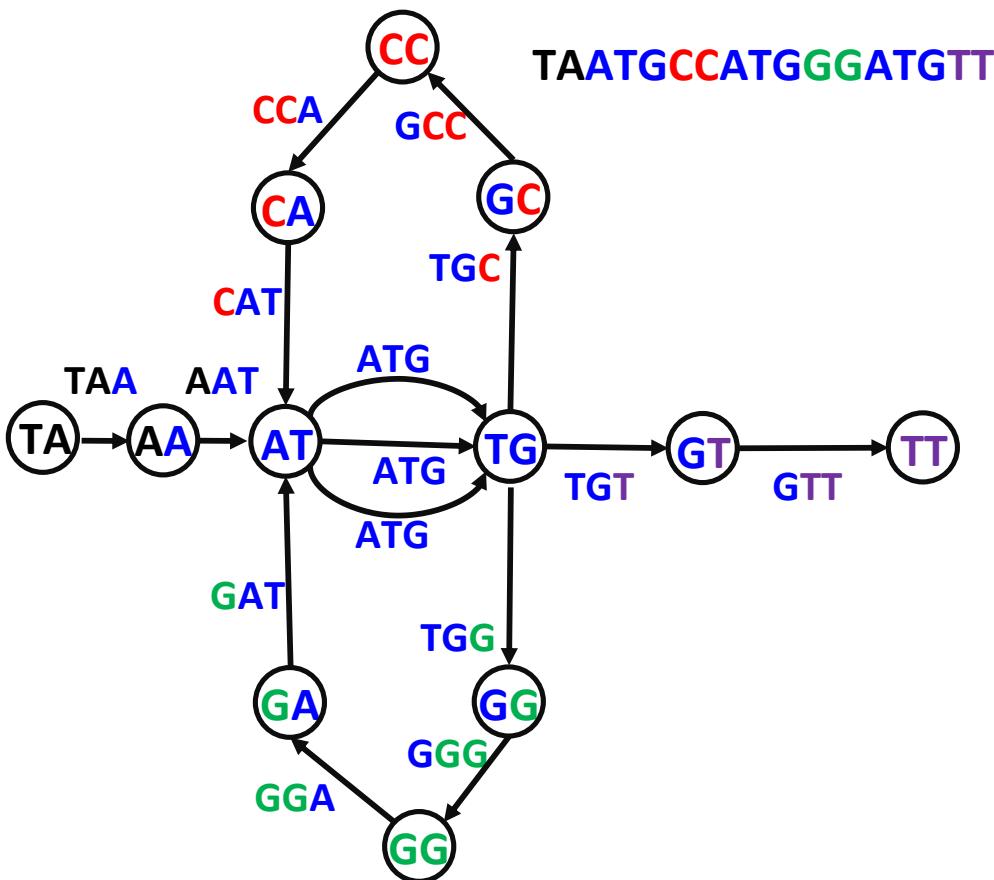
Zalepimo identično obeležene čvorove



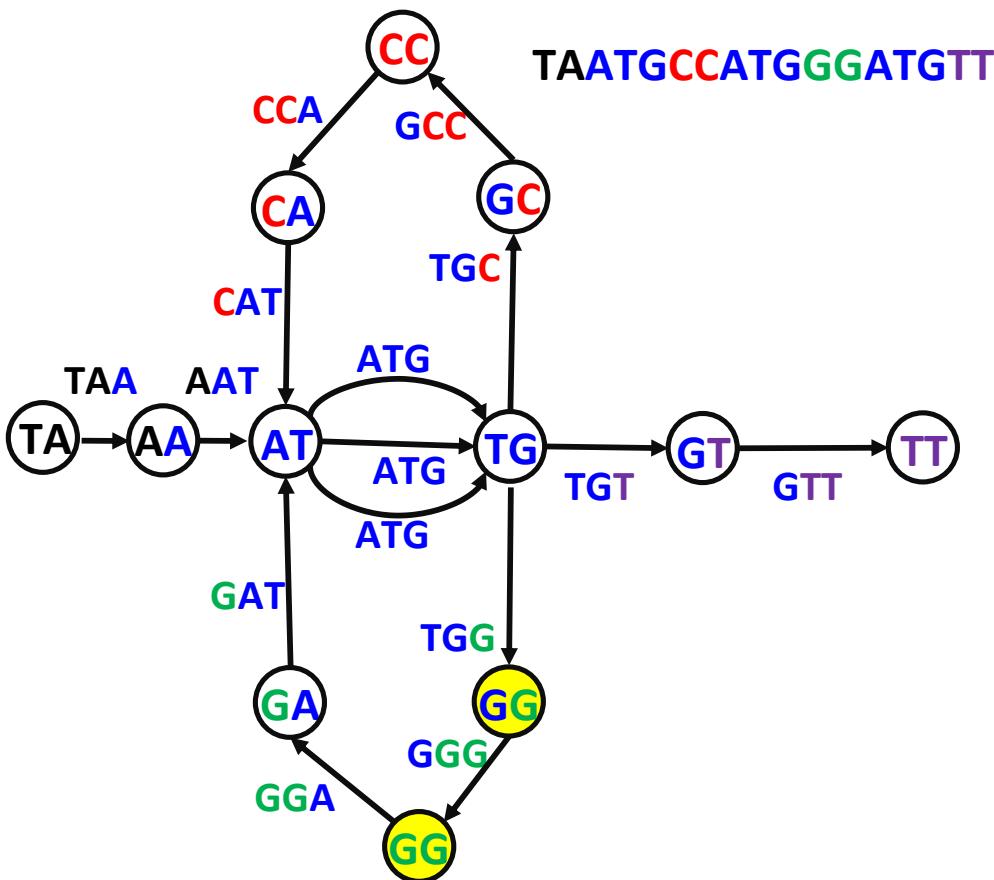
Zalepimo identično obeležene čvorove



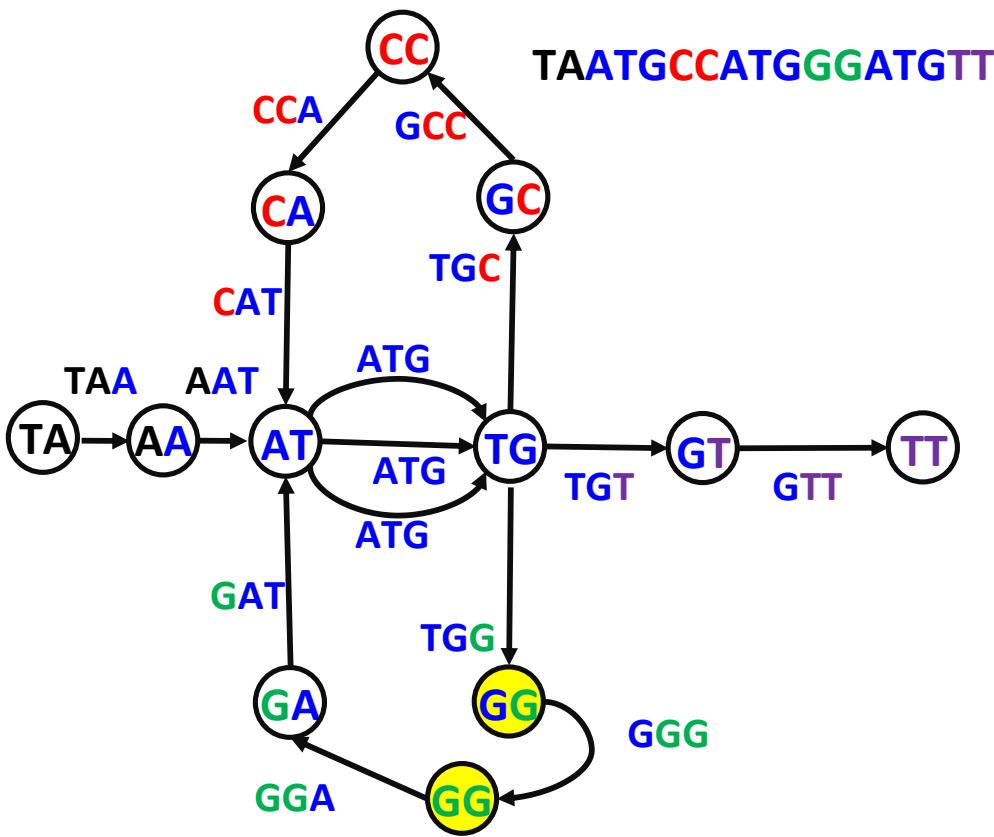
Zalepimo identično obeležene čvorove



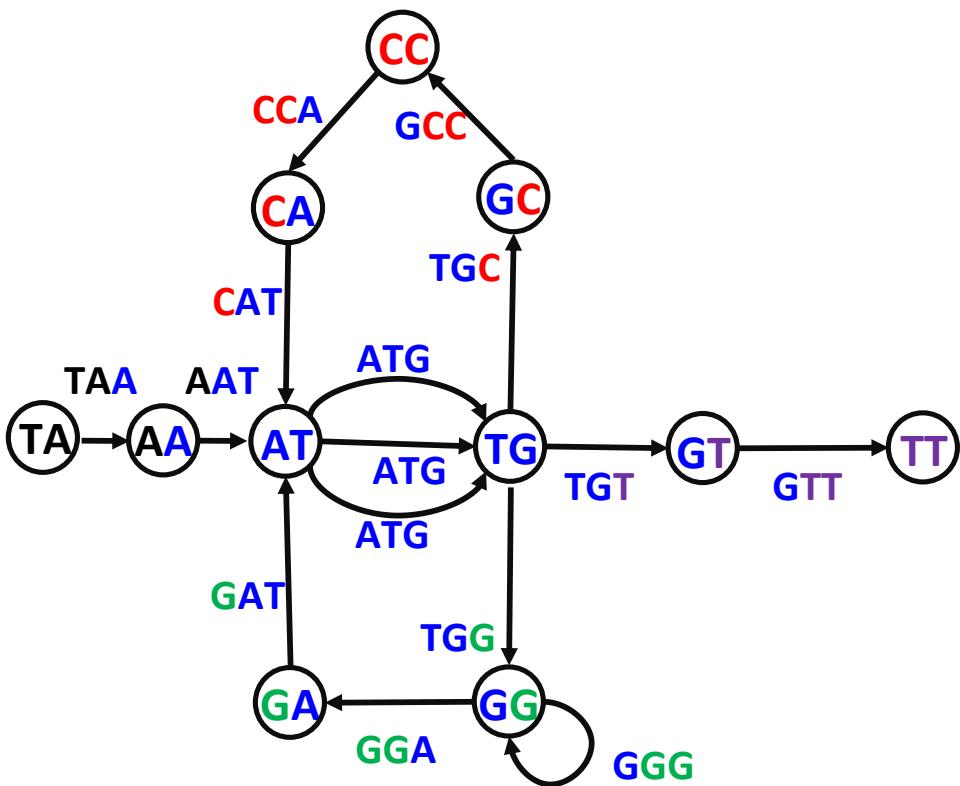
Zalepimo identično obeležene čvorove



Zalepimo identično obeležene čvorove



Isti De Brujinov graf:
 $DeBruin(Genome) =$
 $DeBruin(Genome\ Composition)$



Konstrukcija De Brujinovog grafa

De Brujinov graf na osnovu kolekcije k -grama:

- Svaka grana je označena jednim k -gramom
- Svaki čvor je označen prefiksom/sufiksom izlazne/ulazne grane
- Zalepljeni su svi čvorovi sa identičnim oznakama.

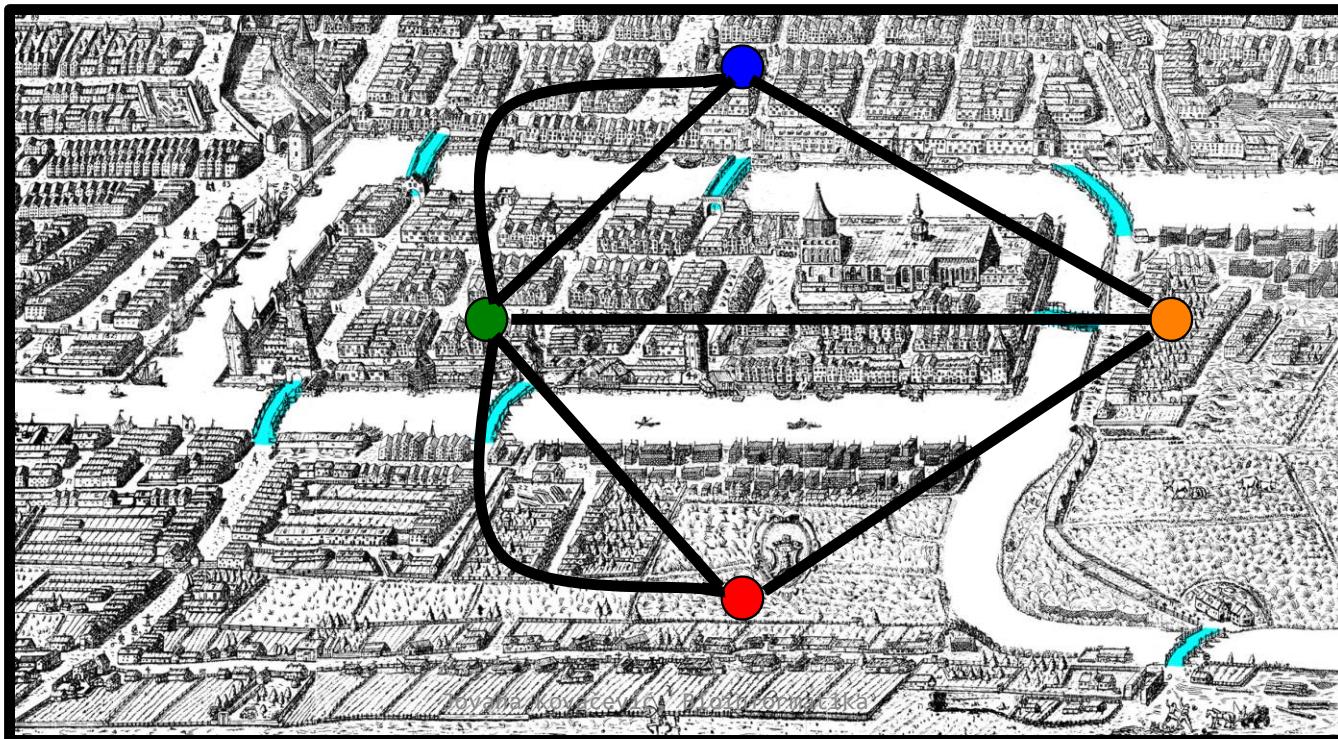
Pregled

- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- Problem rekonstrukcije niske
- Rekonstrukcija niske kao problem Hamiltonove putanje
- Rekonstrukcija niske kao problem Ojlerove putanje
- De Bruijinovi grafovi
- **Ojlerova teorema**
- Spajanje parova očitavanja
- U realnosti

Problem Ojlerovog ciklusa

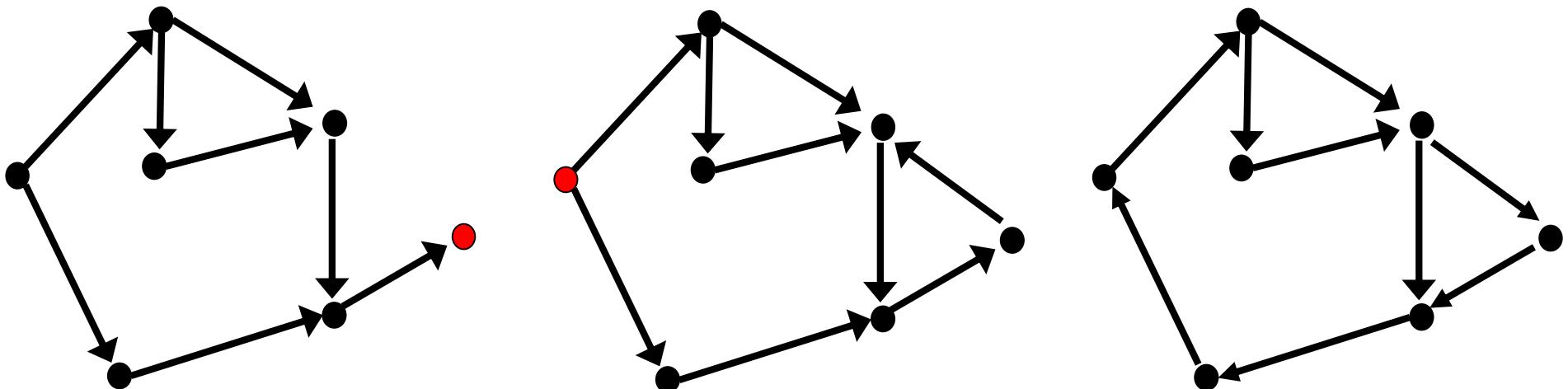
Problem Ojlerovog ciklusa. Pronaći Ojlerov ciklus u grafu.

- **Ulaz.** Graf.
- **Izlaz.** Ciklus koja posećuje svaku granu u grafu tačno jednom.



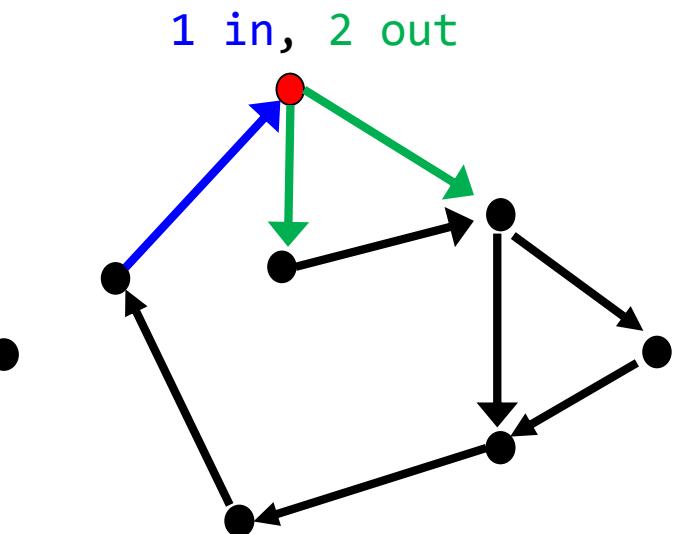
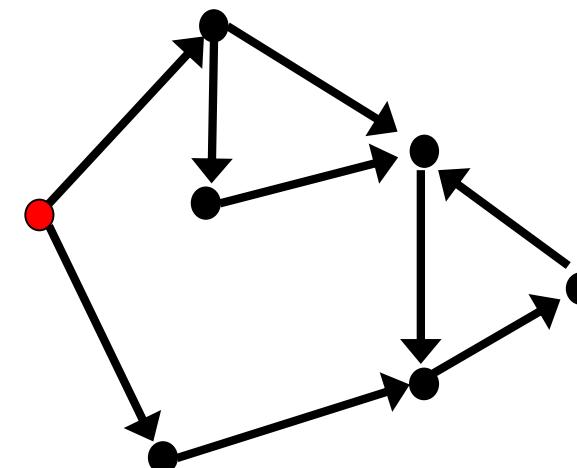
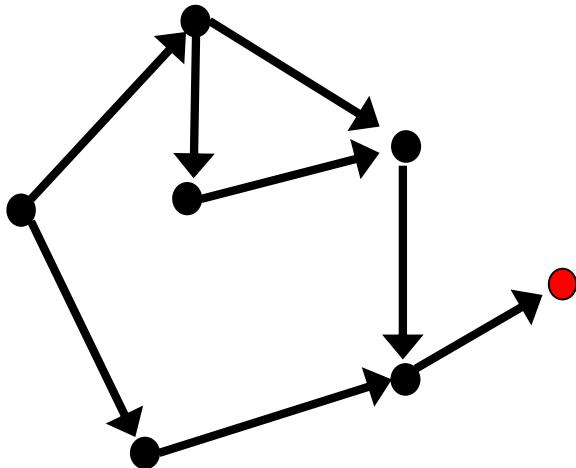
Kažemo da je graf **Ojlerov** ako sadrži Ojlerov ciklus.

Da li je ovaj graf Ojlerov?



Kažemo da je graf Ojlerov ako sadrži Ojlerov ciklus.

Da li je ovaj graf Ojlerov?



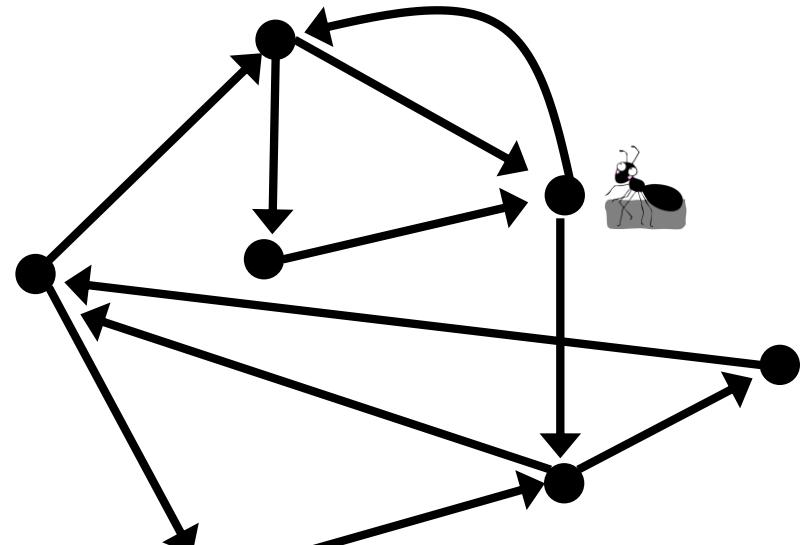
Kažemo da je graf balansiran ako za svaki čvor važi
indegree = outdegree

Ojlerova teorema

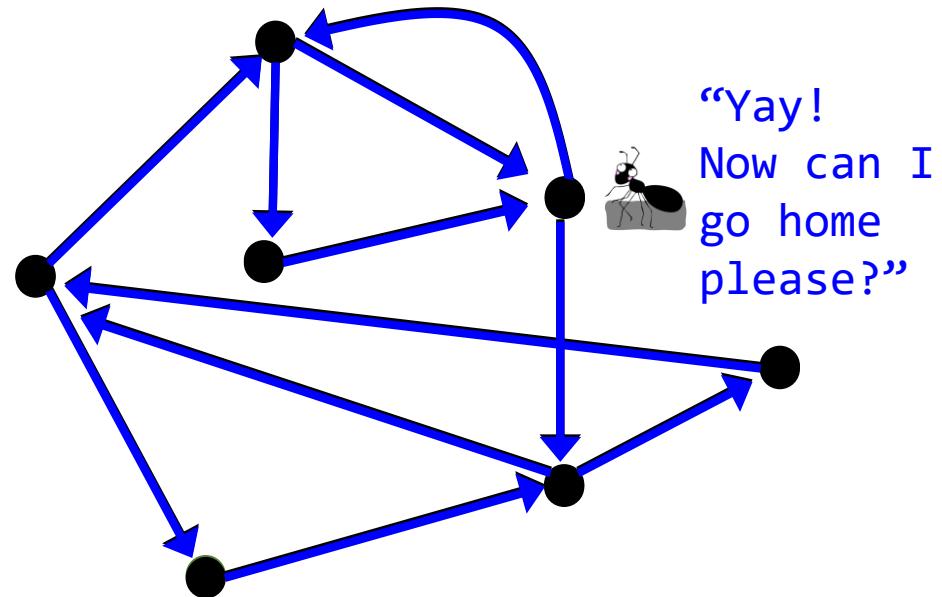
- Svaki Ojlerov graf je balansiran
- Svaki povezan i balansiran graf je Ojlerov
 - Kažemo da je graf povezan ako za mrežu koja dva čvora postoji putanja koja ih povezuje.

Kako bi mrav dokazao Ojlerovu teoremu?

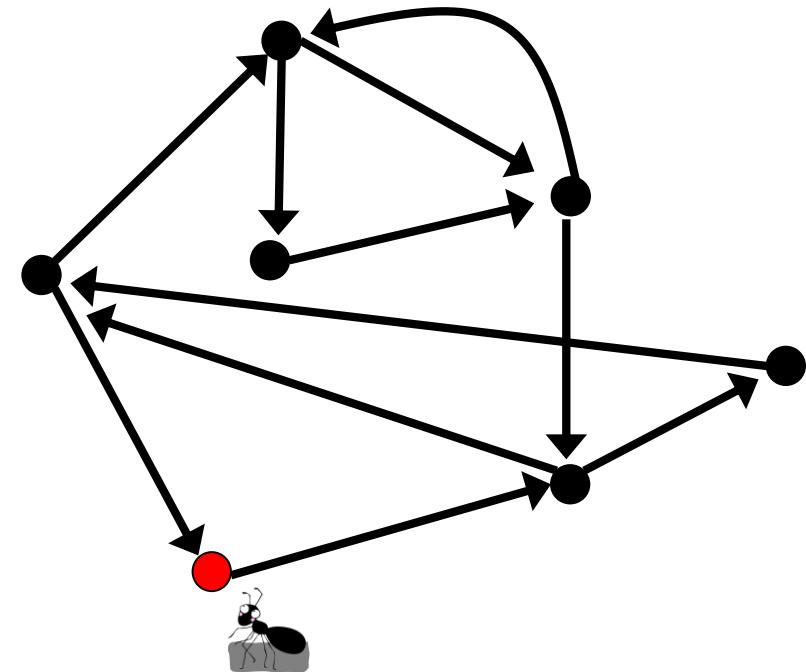
Mrav na slučajan način bira grane kojima će se kretati u grafu. **Ne može da obiđe istu granu dvaput!**



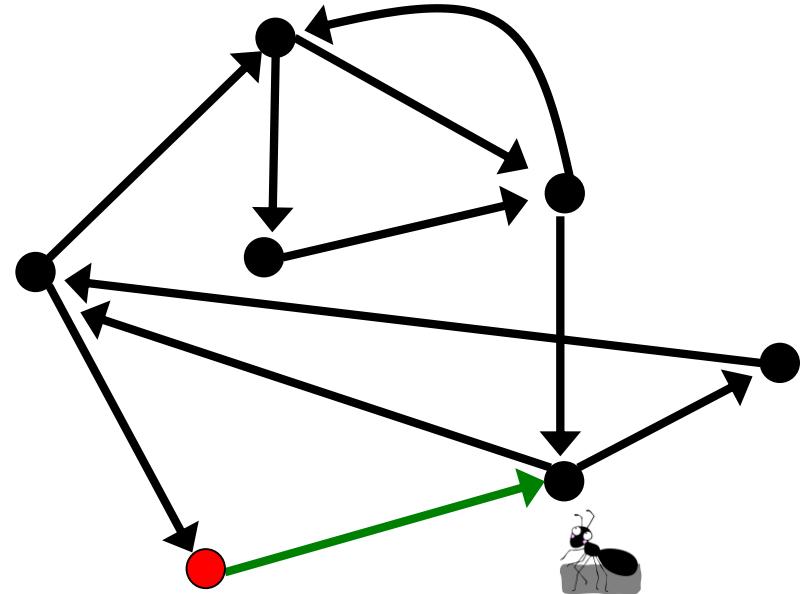
Veoma pametan mrav



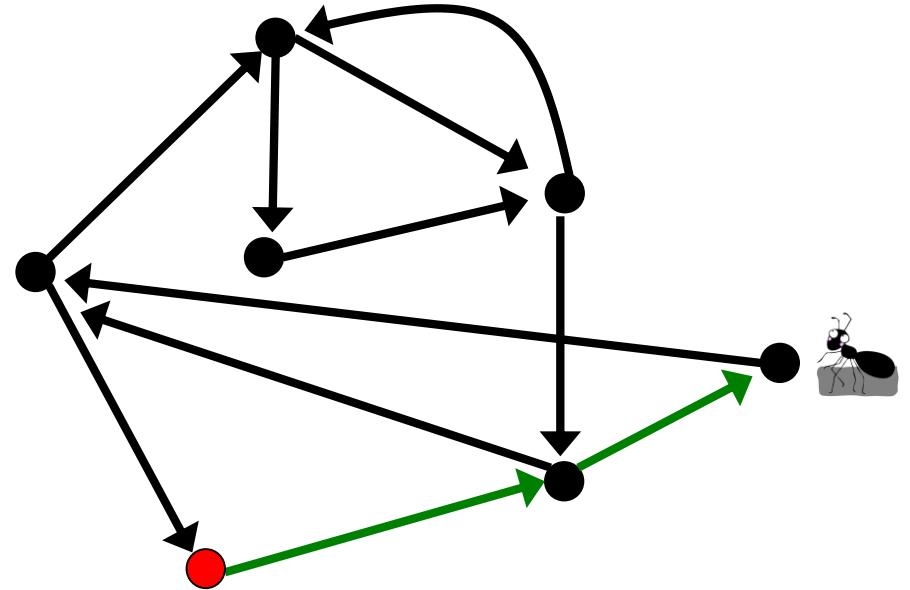
Manje pametan mrav



... obilazi ...

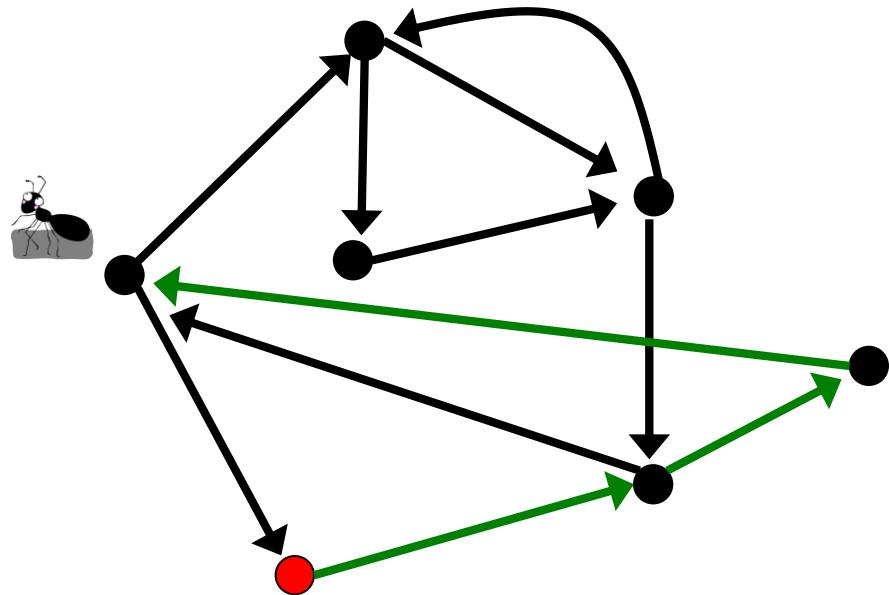


... obilazi ...

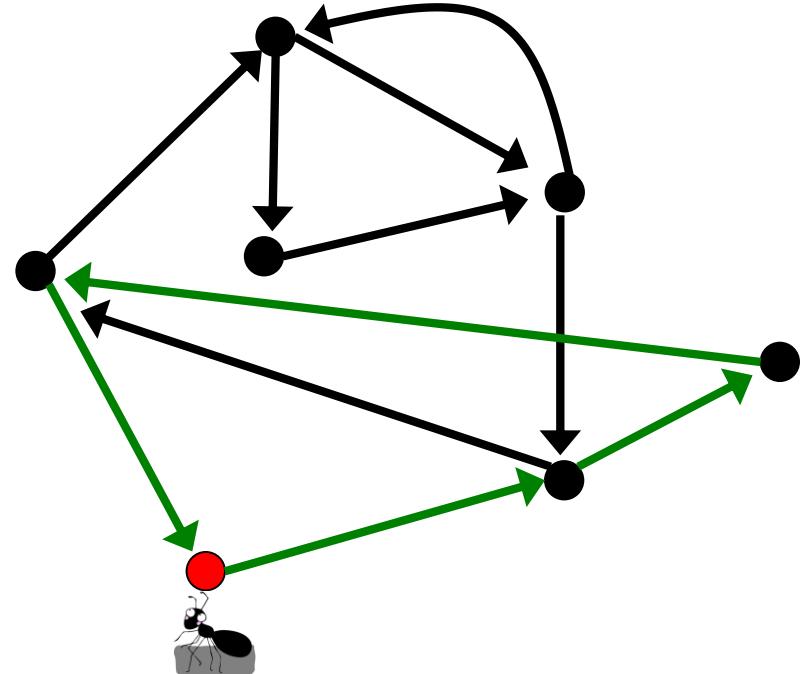


... obilazi ...

Da li može da se zaglavi? **U kom čvoru?**

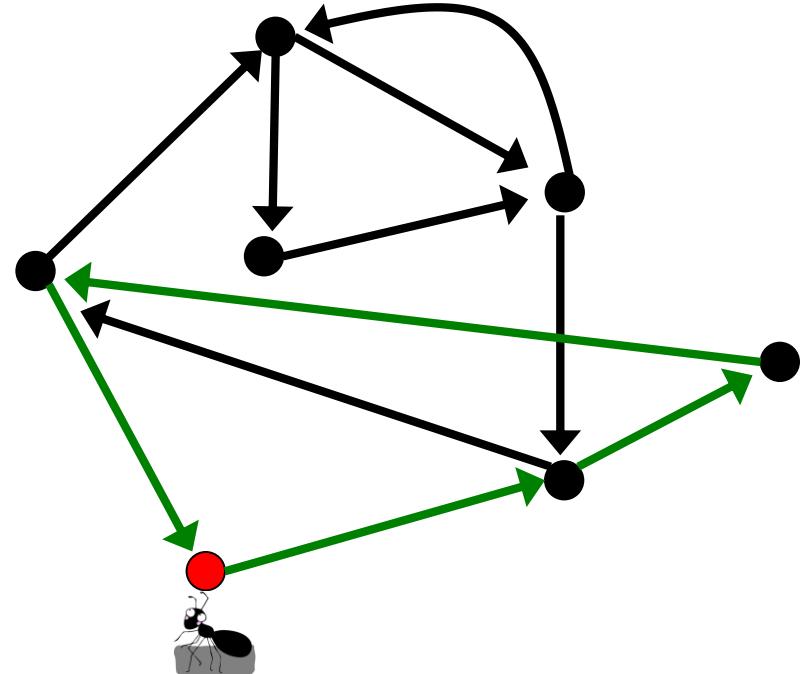


Mrav može da se zaglavi samo u čvoru iz kog je počeo obilazak



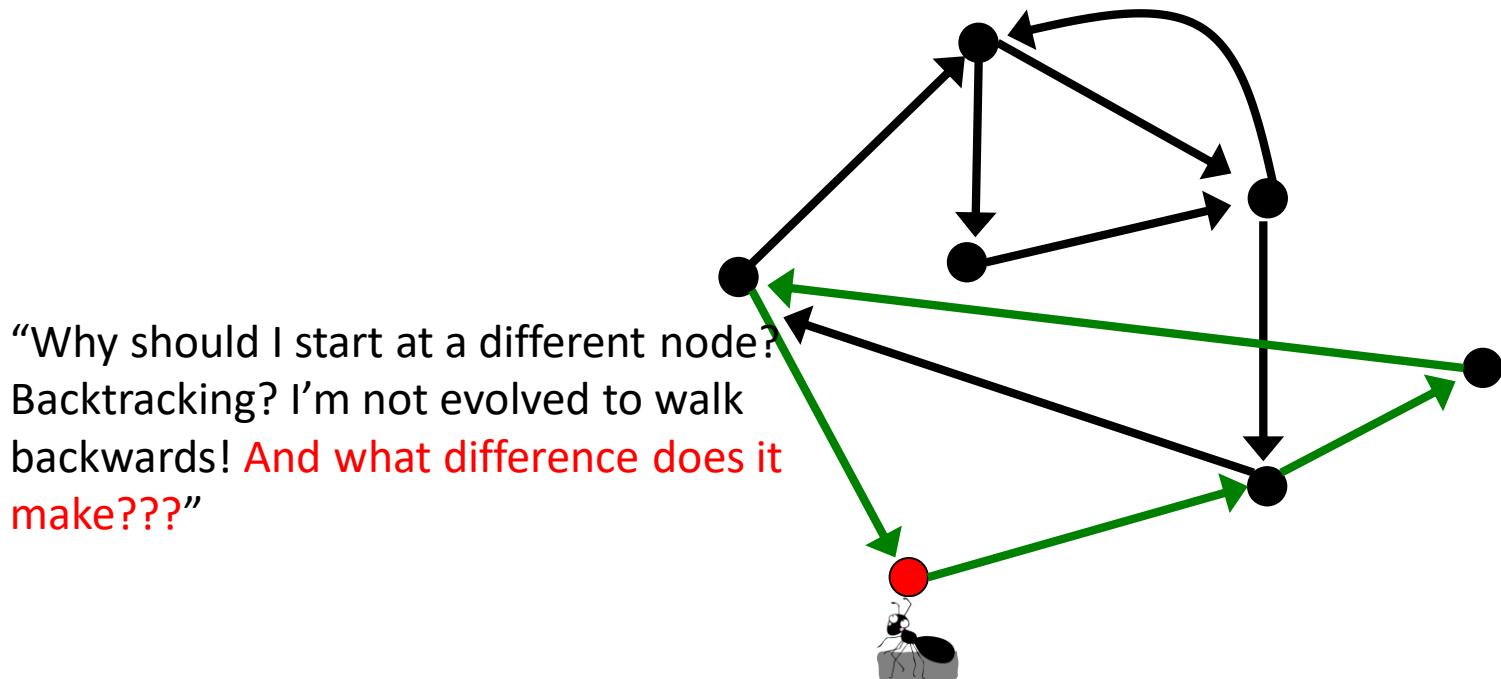
Mrav je kreirao ciklus ali još nije dokazao Ojlerovu teoremu

Konstruisani ciklus nije Ojlerov. **Možemo li da ga uvećamo?**



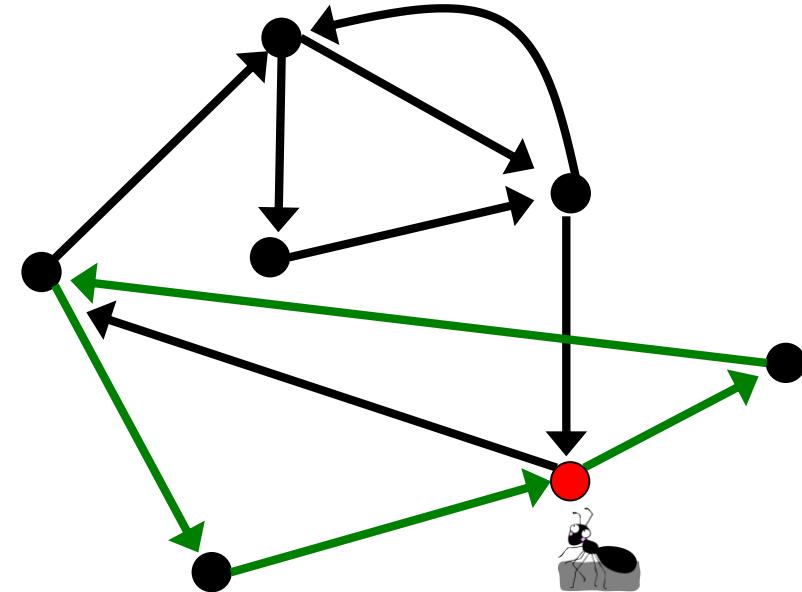
Hajde da započnemo obilazak u nekom drugom čvoru iz zelenog ciklusa

U kom? U onom koji ima neposećene grane.



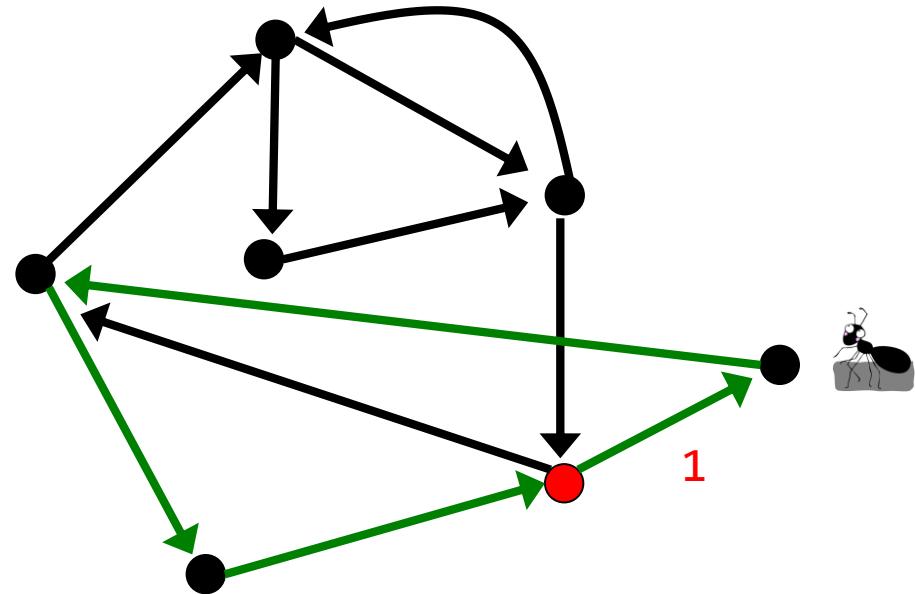
Nove instrukcije za mrava:

Počni od **čvora** koji ima neposećenu granu, obidi već konstruisani zeleni ciklus i vrati se u početni čvor



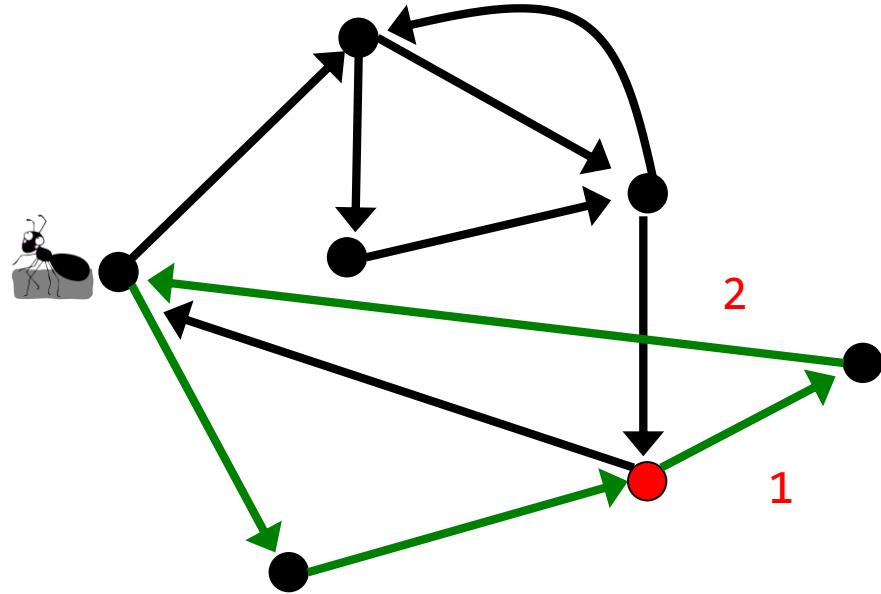
Mrav obilazi prethodno konstruisani ciklus

Počni od čvora koji ima neposećenu granu, obidi već konstruisani zeleni ciklus i vratи se u početni čvor



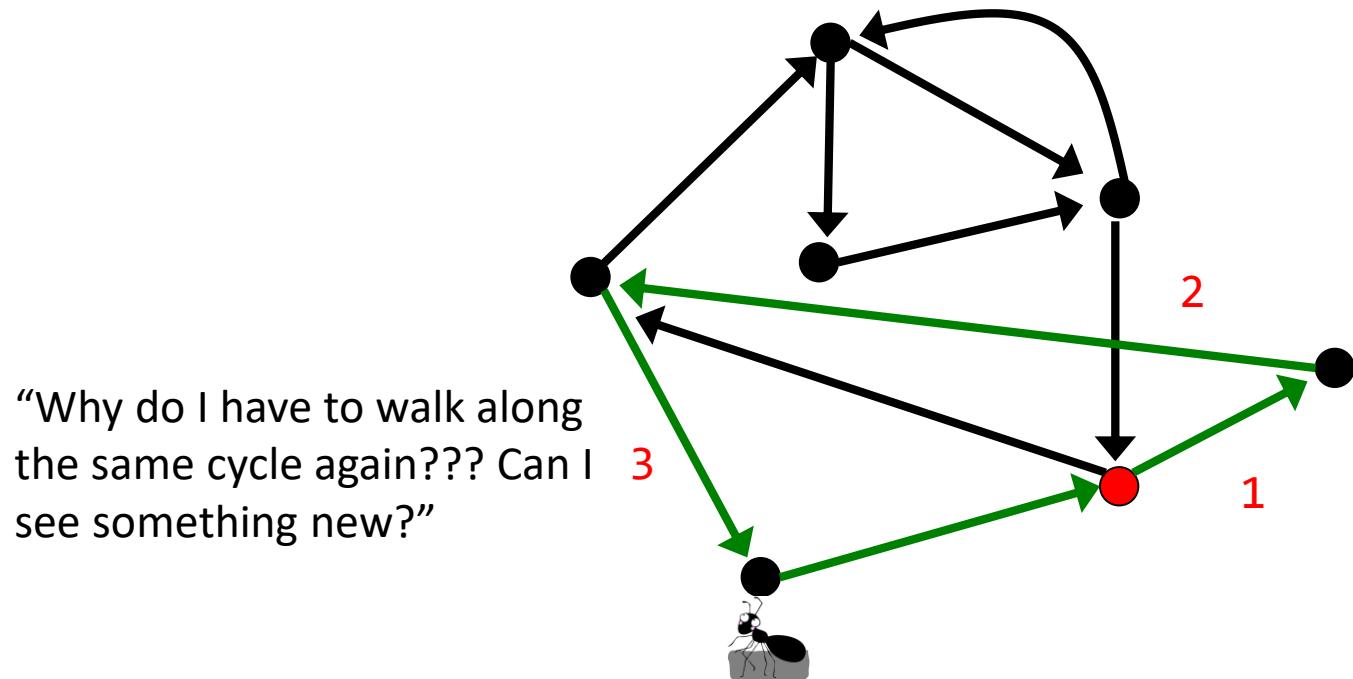
Mrav obilazi prethodno konstruisani ciklus

Počni od čvora koji ima neposećenu granu, obidi već konstruisani zeleni ciklus i vratи se u početni čvor



Mrav obilazi prethodno konstruisani ciklus

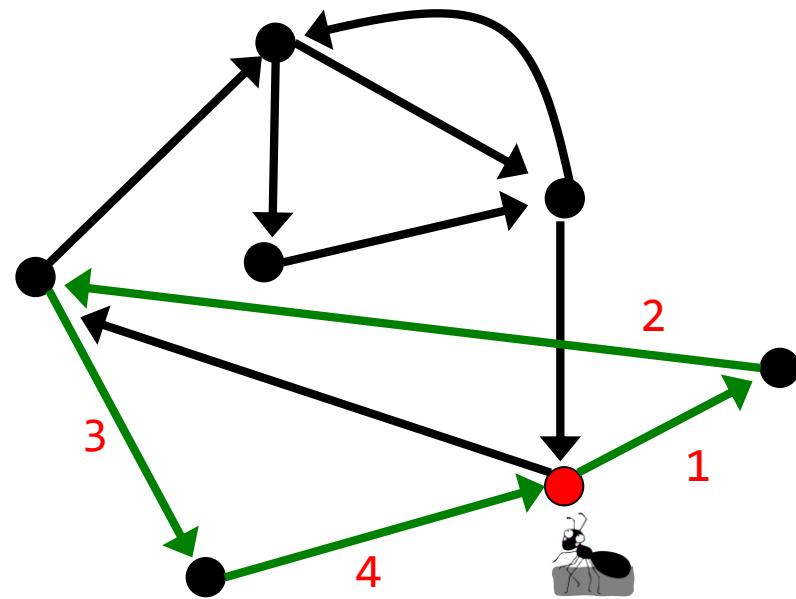
Počni od čvora koji ima neposećenu granu, obidi već konstruisani zeleni ciklus i vratи se u početni čvor



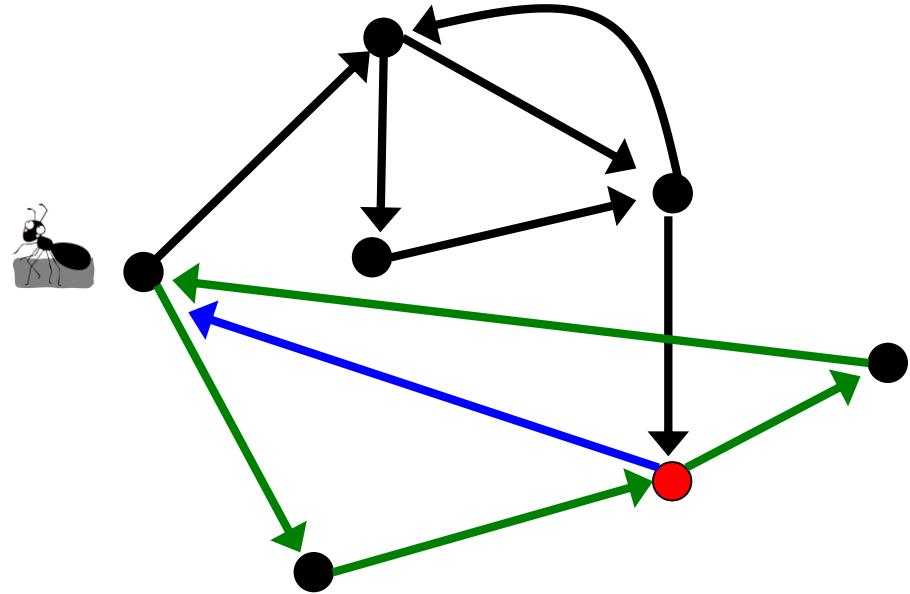
Mrav se vratio nazad ali može da nastavi da obilazi!

Počni od čvora koji ima neposećenu granu, obidi već konstruisani zeleni ciklus i vrati se u početni čvor

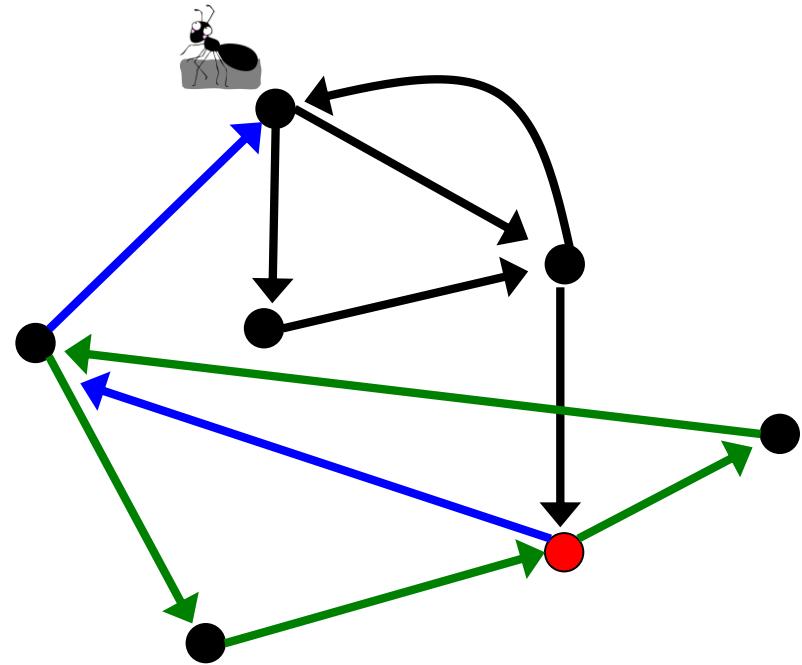
Nakon završenog ciklusa, nastavi obilazak tako što ćeš posetiti neku granu koja nije posećena ranije. Ako ih ima više, odaberi jednu na slučajan način.



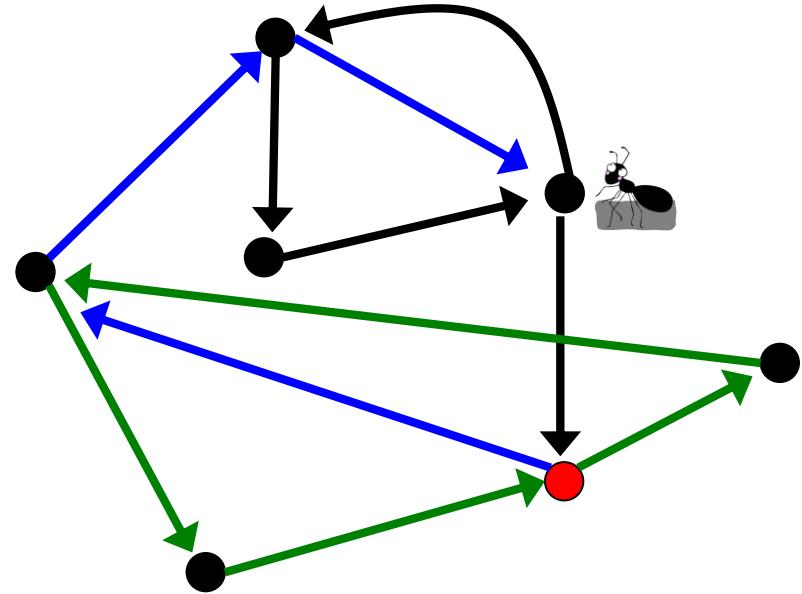
Uvećavamo prethodno konstruisani ciklus



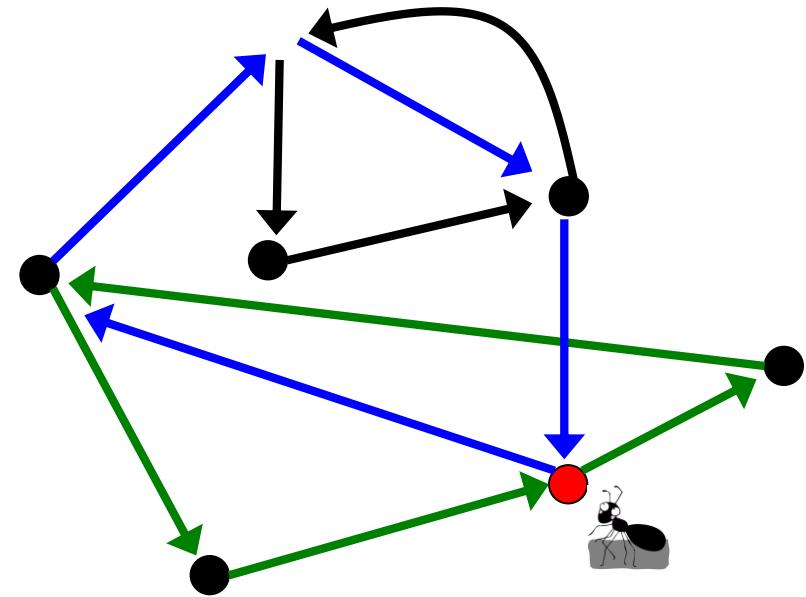
Uvećavamo prethodno konstruisani ciklus



Uvećavamo prethodno konstruisani ciklus



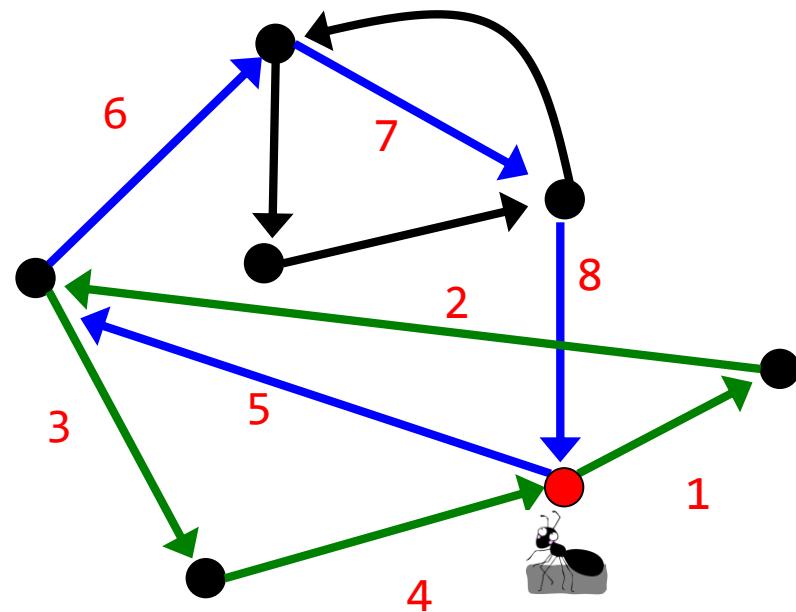
Uvećavamo prethodno konstruisani ciklus



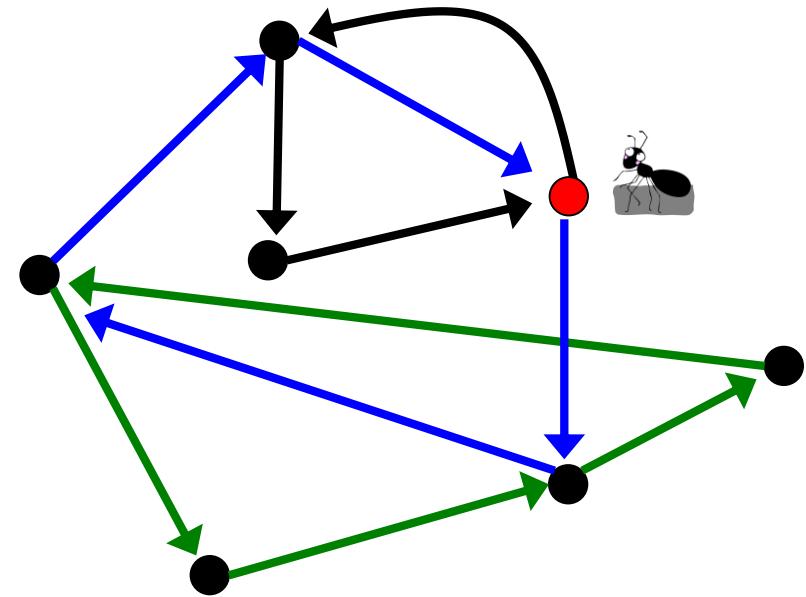
Mrav se ponovo zaglavio!

Konstruisani zeleno-plavi ciklus i dalje nije Ojlerov. Da li možemo da ga uvećamo?

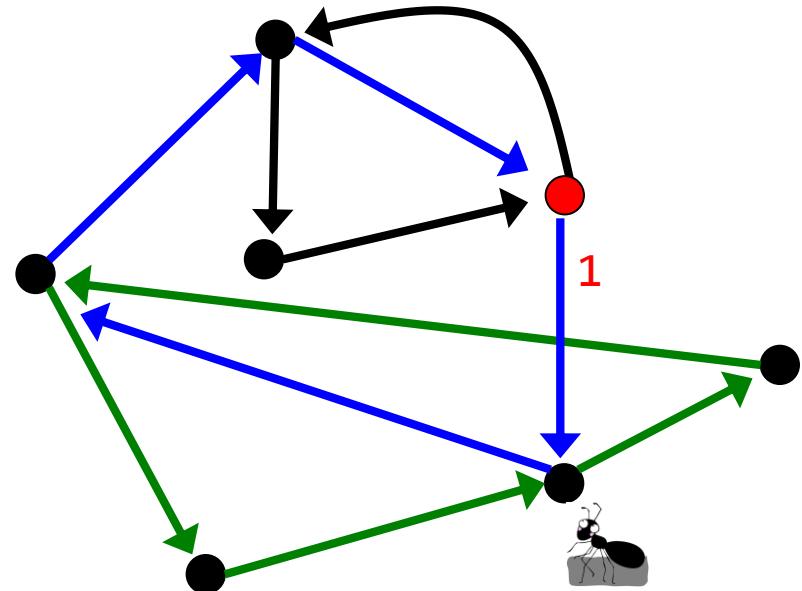
Mrav treba da obide konstruisani ciklus počev od drugog čvora. Od kog?



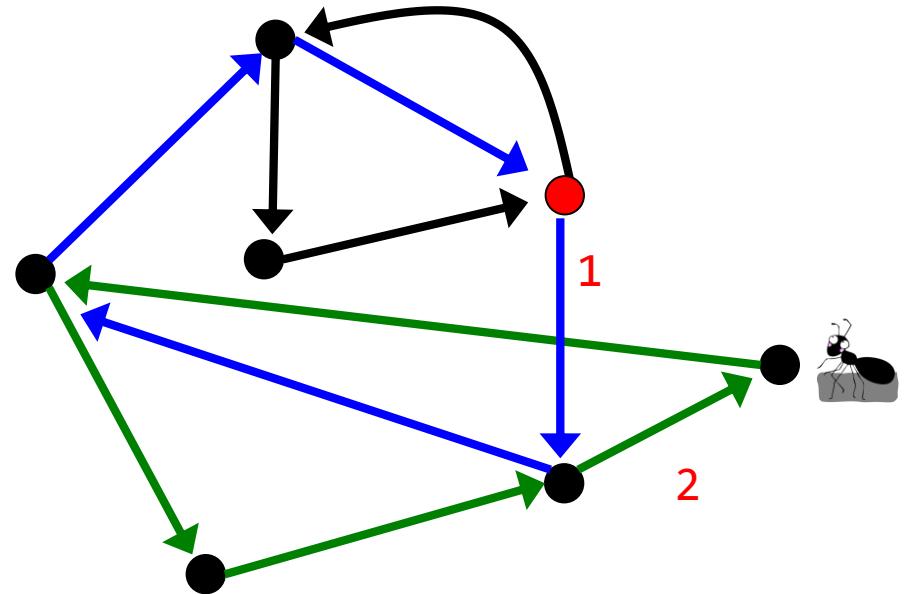
Počinjemo od novog čvora, ponovo...



Obilazimo prethodno konstruisan zeleno-plavi ciklus

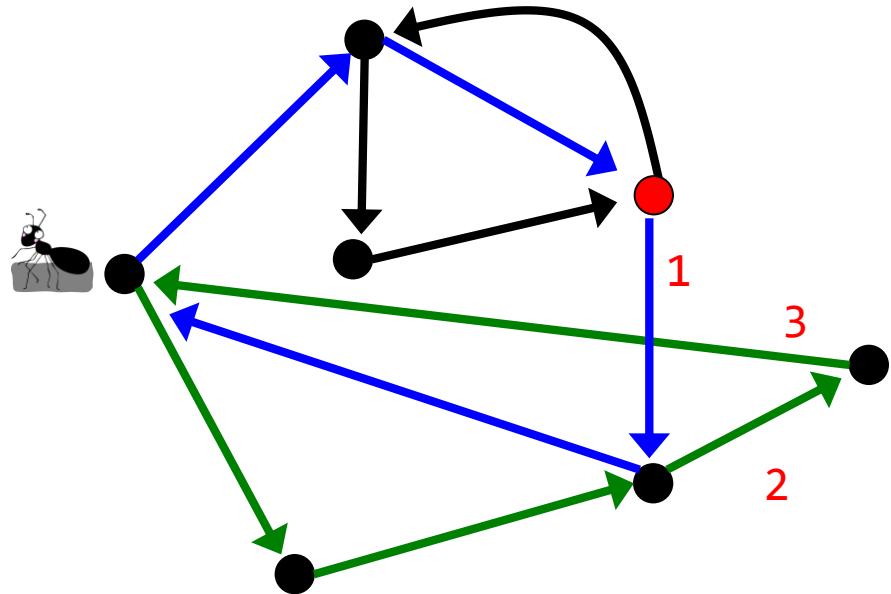


Obilazimo prethodno konstruisan zeleno-plavi ciklus

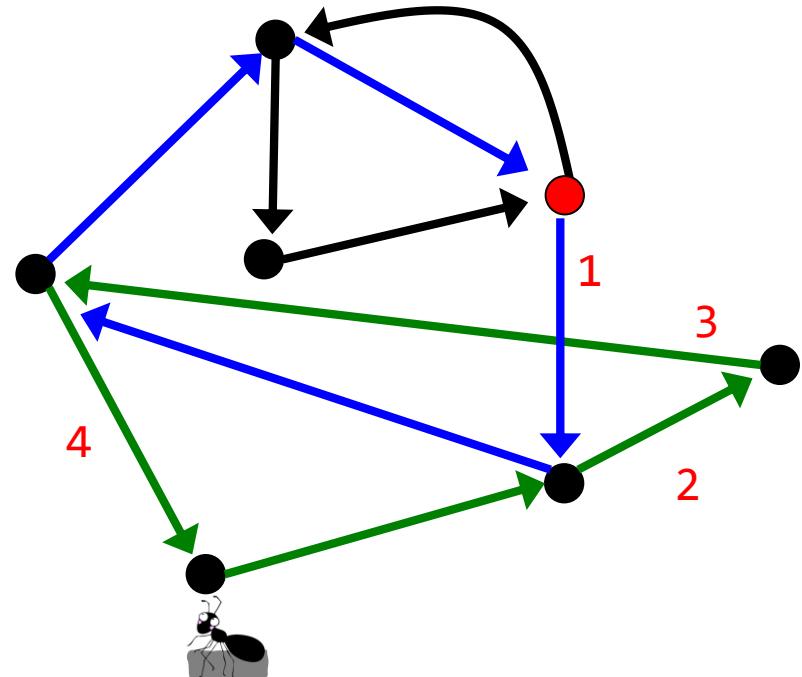


Obilazimo prethodno konstruisan zeleno-plavi ciklus

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

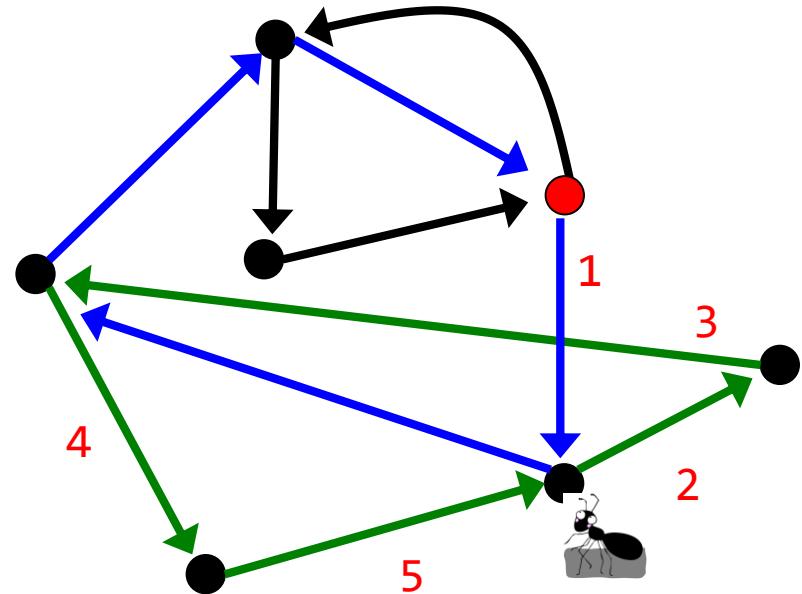


Obilazimo prethodno konstruisan zeleno-plavi ciklus

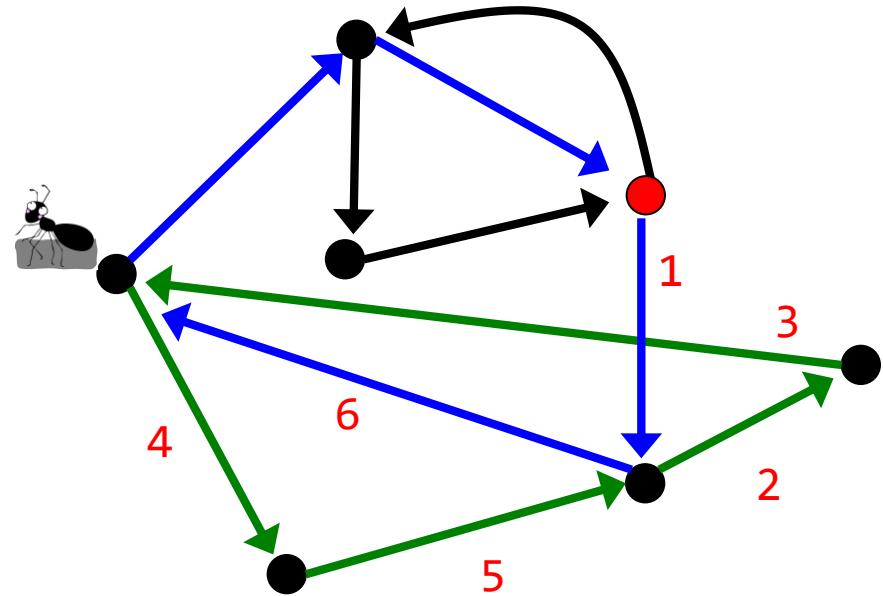


“These instructions are stupid...”

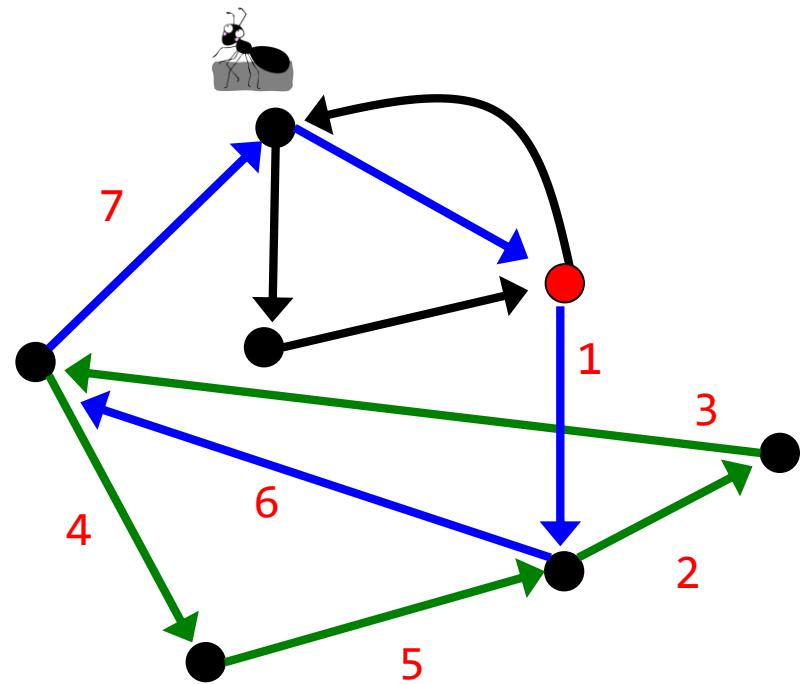
Obilazimo prethodno konstruisan zeleno-plavi ciklus



Obilazimo prethodno konstruisan zeleno-plavi ciklus

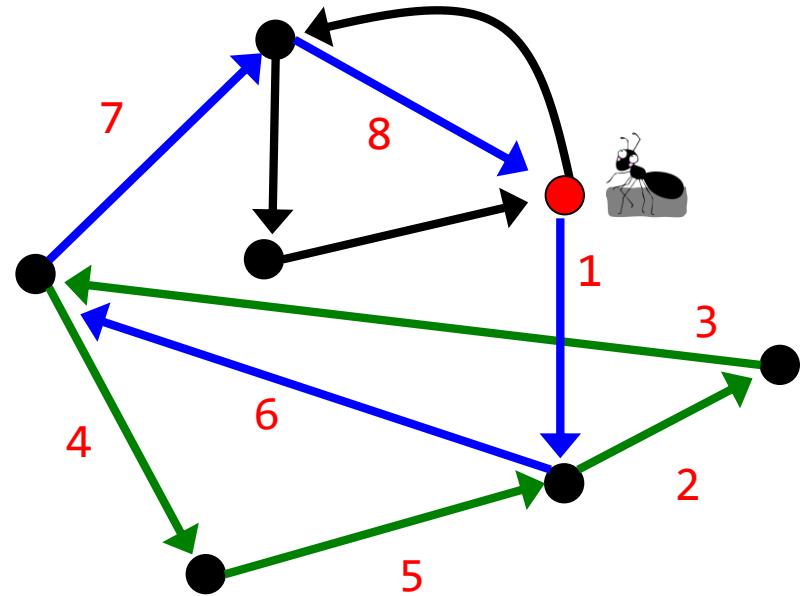


Obilazimo prethodno konstruisan zeleno-plavi ciklus

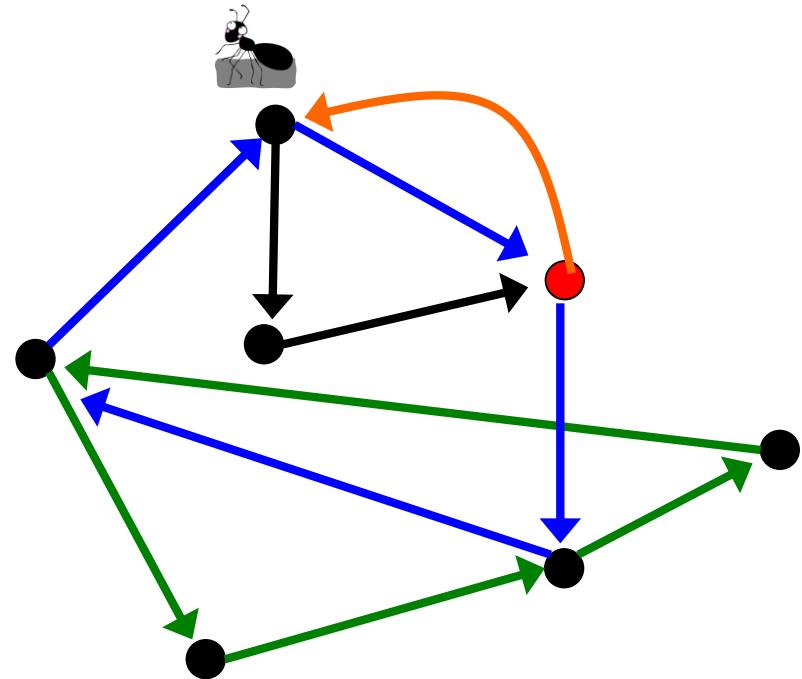


Mrav se vratio nazad ali može da nastavi da obilazi!

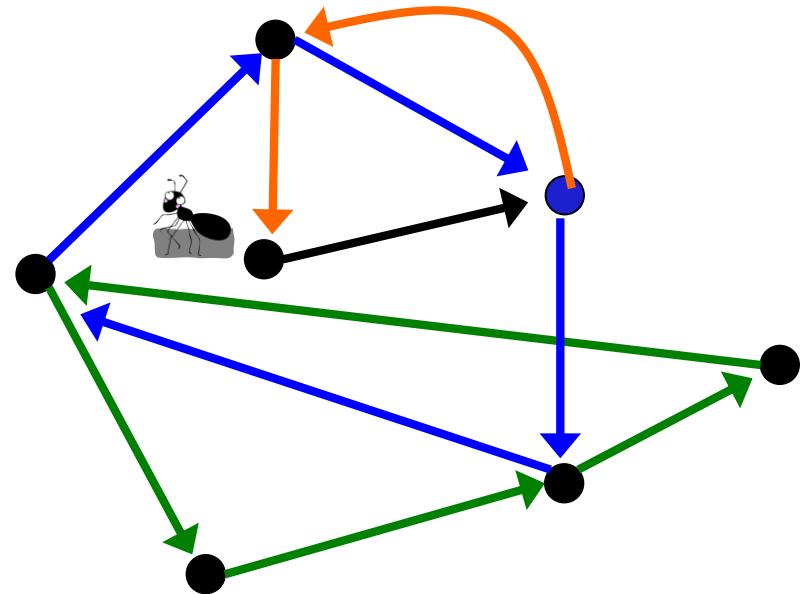
"Hmm, maybe these instructions were not that stupid..."



Uvećavamo zeleno-plavi ciklus



Uvećavamo zeleno-plavi ciklus



Ojlerova teorema je dokazana

EulerianCycle(BalancedGraph)

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

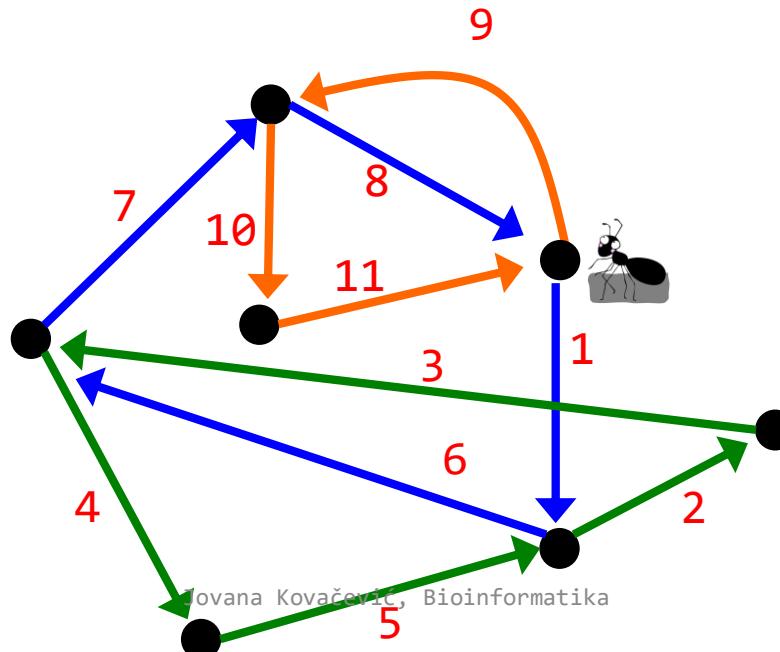
while *Cycle* is not Eulerian

 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*

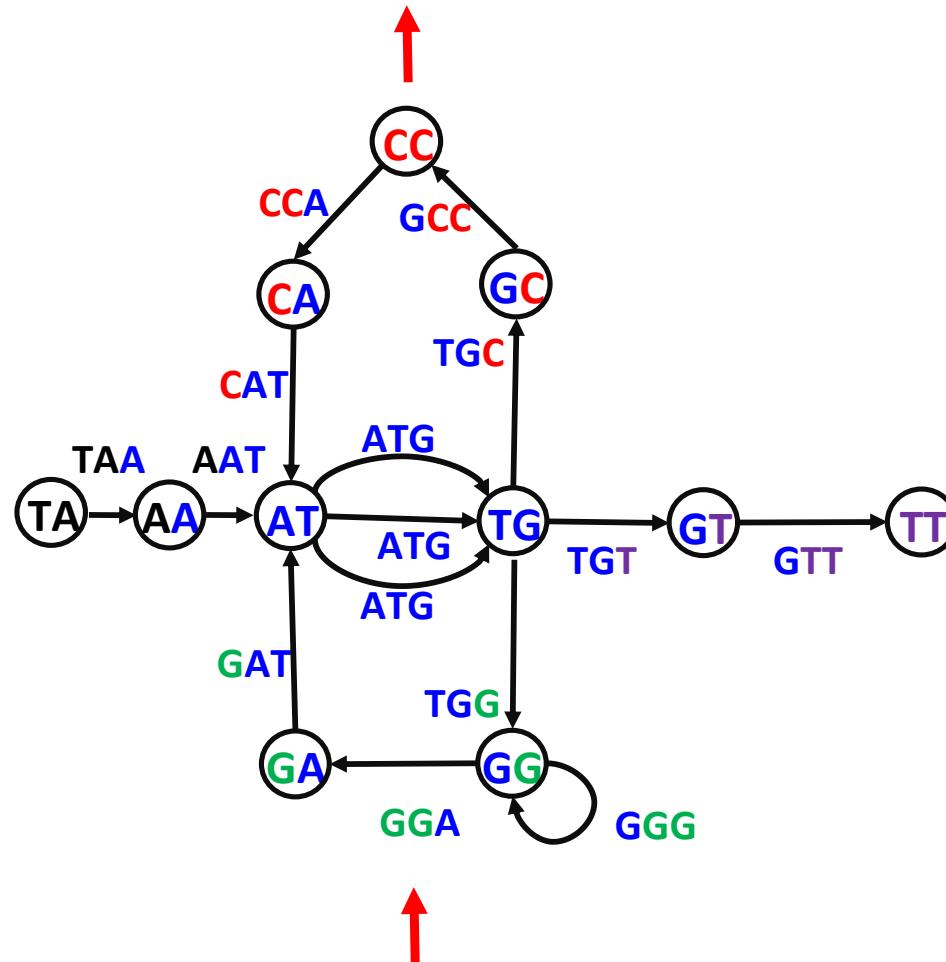


Pregled

- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- Problem rekonstrukcije niske
- Rekonstrukcija niske kao problem Hamiltonove putanje
- Rekonstrukcija niske kao problem Ojlerove putanje
- Slični problemi sa različitim sudbinama?
- De Bruijnovi grafovi
- Ojlerova teorema
- **Sastavljanje parova očitavanja**
- De Bruijn Graphs Face Harsh Realities of Assembly

Od očitavanja do De Brujinovog grafa do genoma

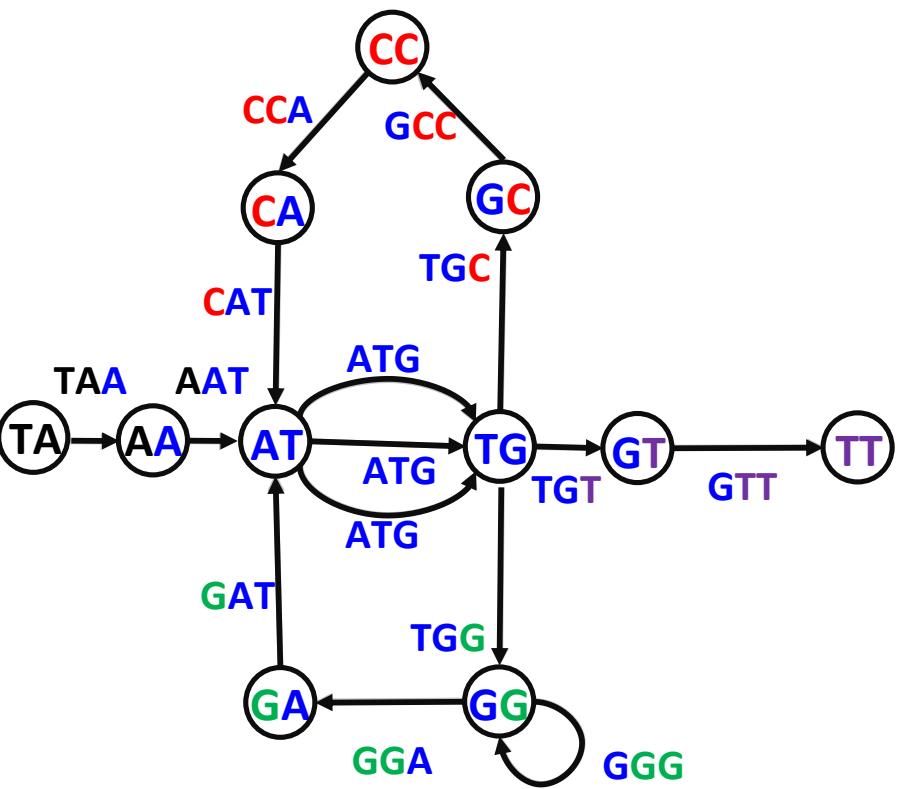
TA~~ATGCCATGGGATGTT~~



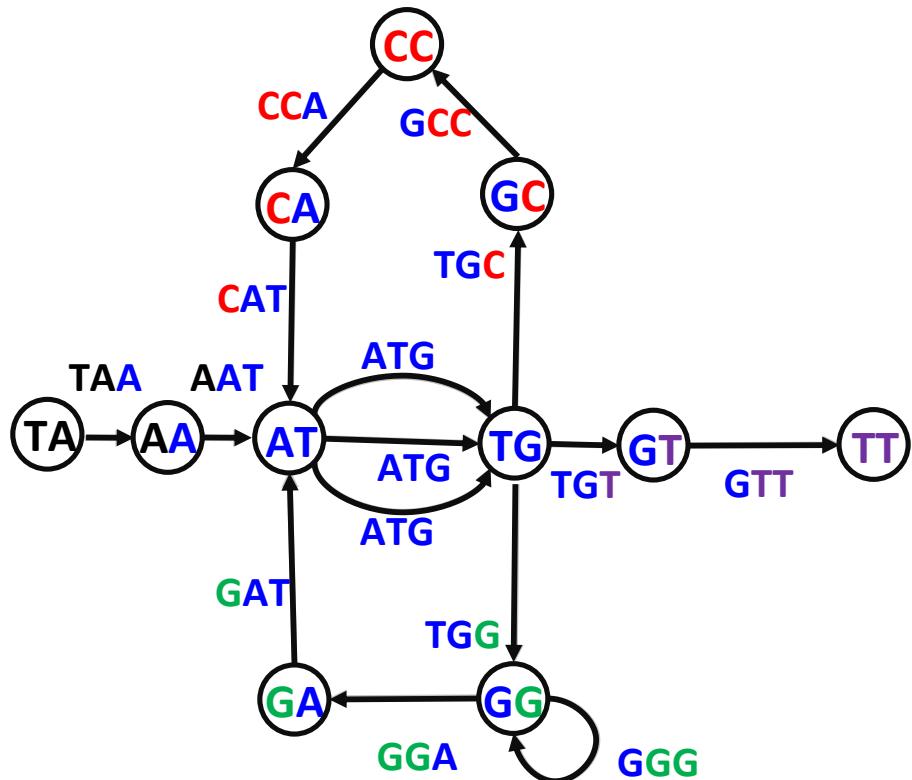
AAT ATG ATG ATG CAT CCA GAT GCC GGA GGG GTT TAA TGC TGG TGT

Graf može imati više Ojlerovih putanja

TAATGCCATGGGATGTT

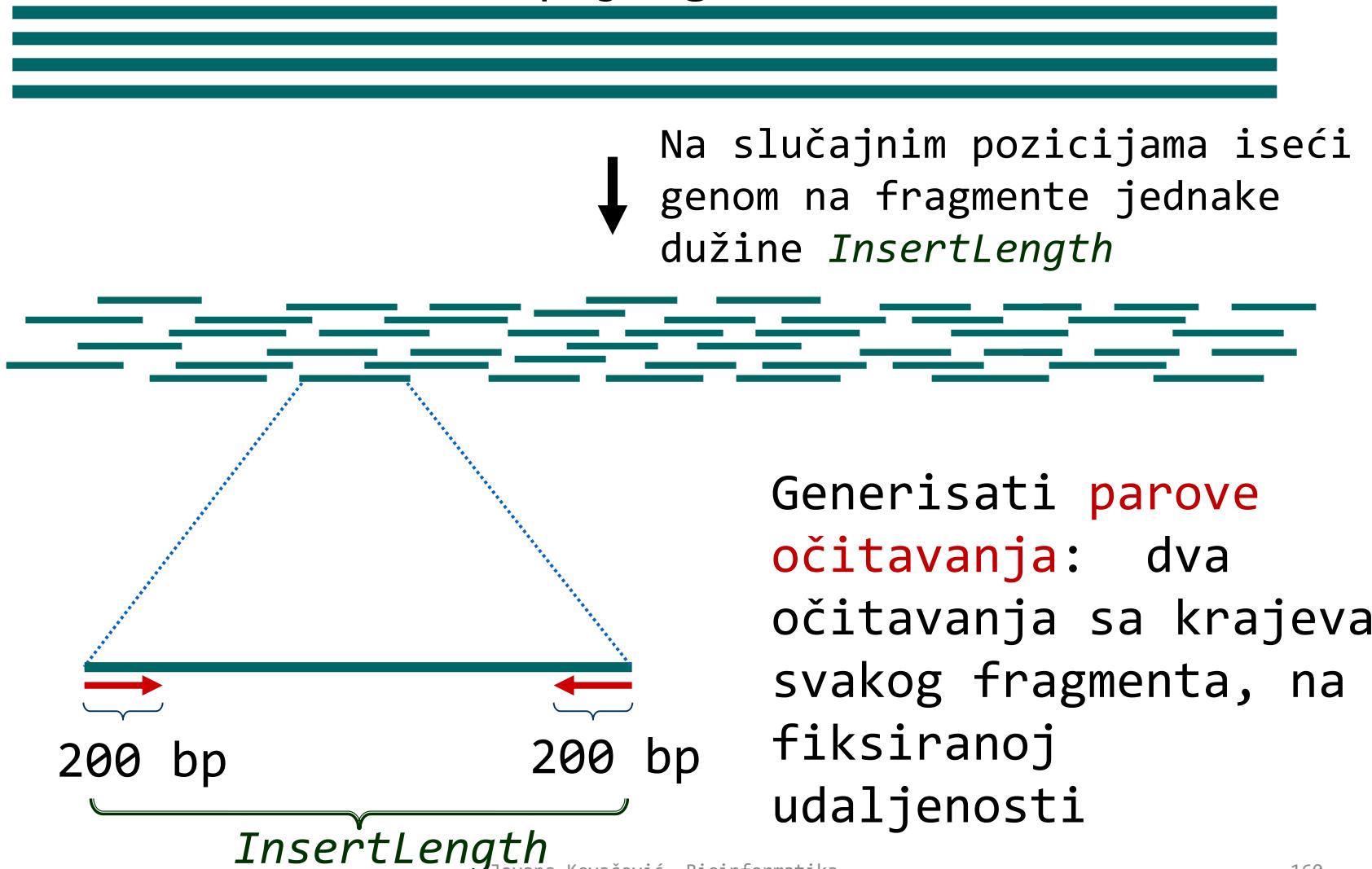


TAATGGGATGCCATGTT



DNK sekvencioniranje sa parovima očitavanja

Više identičnih kopija genoma



Od k -grama do uparenih k -grama



Pod uparenim k -gramom podrazumevamo par k -grama na fiksiranom rastojanju d u genomu. Na primer, **TCA** i **TCC** na rastojanju $d=11$ čine jedan upareni k -gram.

Šta je upareni k -gramske sastav PairedComposition(**TAATGCCATGGGATGTT**)?

TAA **GCC**

upareni 3-gram

Šta je upareni k -gramske sastav PairedComposition(**TAATGCCATGGGATGTT**)?

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

Predstavimo upareni 3-gram **TAA** **GCC** na sledeći način:

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

Leksikografski poredak kolekcije PairedComposition

Problem rekonstrukcije niske na osnovu parova očitavanja

Problem rekonstrukcije niske na osnovu parova očitavanja .
Rekonstruisati nisku na osnovu njenih uparenih k -grama.

- **Ulaz.** Kolekcija uparenih k -grama.
- **Izlaz.** Niska $Text$ takva da je $PairedComposition(Text)$ jednak kolekciji uparenih k -grama.

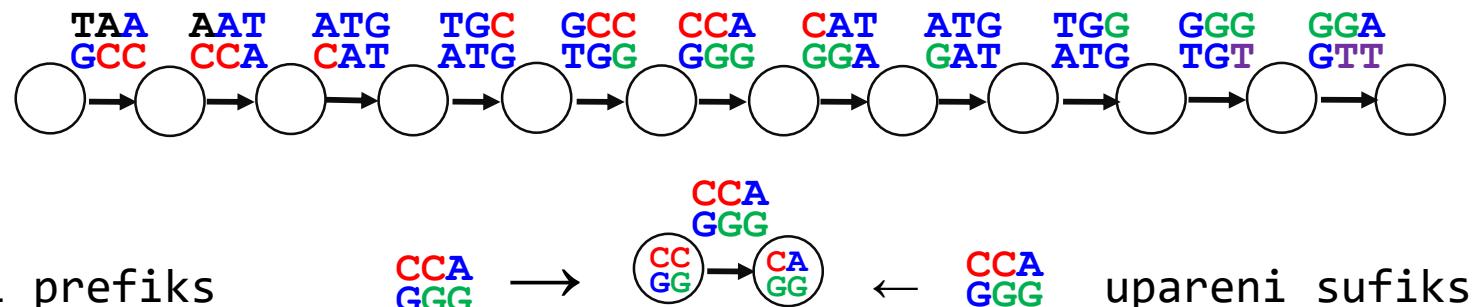
Kako konstruisati upareni De Brujinov graf na osnovu uparenog k -gramskog sastava?



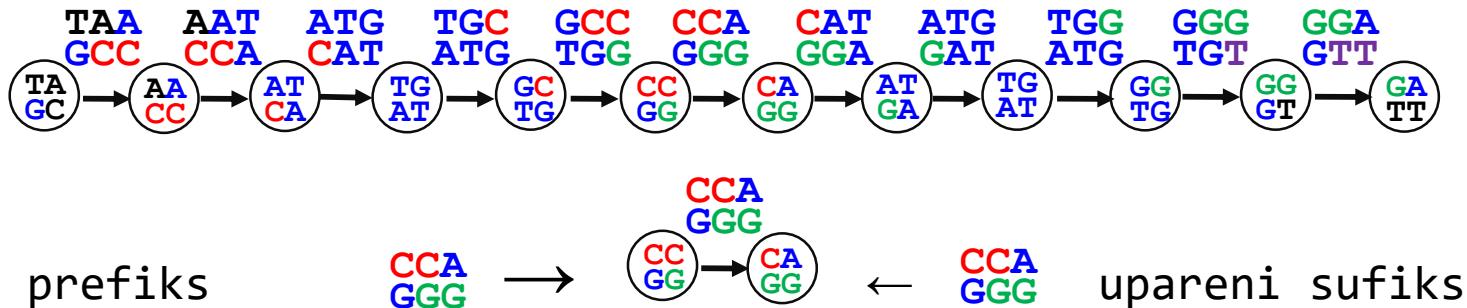
Pretpostavimo da je **dat genom** (niska *Genome*). Posmatrajmo genom kao putanju u grafu obeleženom na osnovu njegovog uparenog k -gramskog sastava

Predstavimo genom **TAATGCCATGGGATGTT** kao putanju

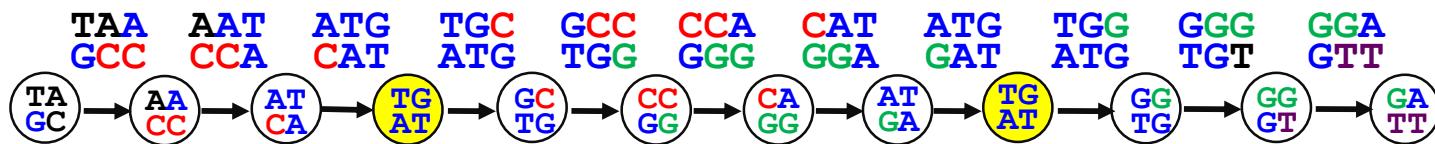
TAA GCC
AAT CCA
ATG CAT
TGC ATG
GCC TGG
CCA GGG
CAT GGA
ATG GAT
TGG ATG
GGG TGT
GGA GTT



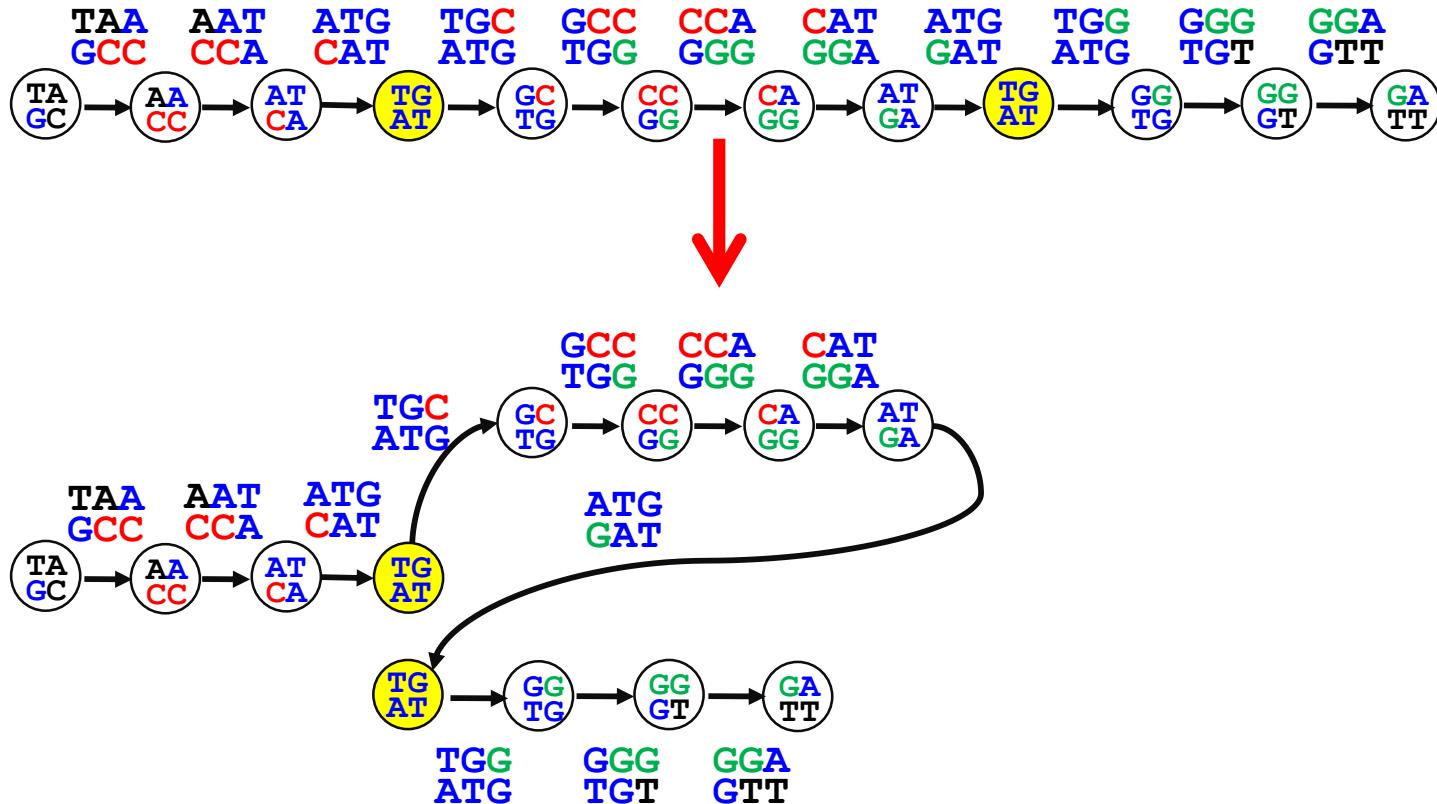
Obeležimo čvorove uparenim prefiksima i sufiksima



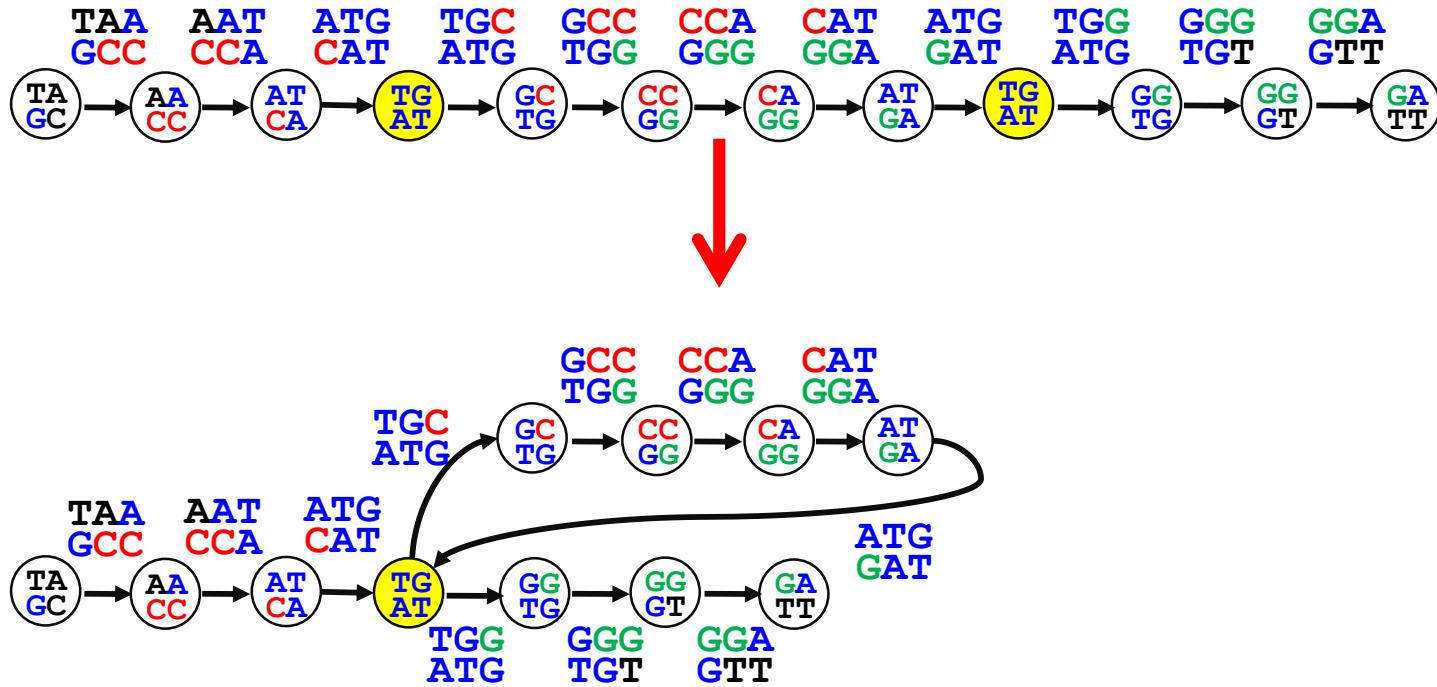
Zalepimo čvorove za identičnim oznakama



Zalepimo čvorove za identičnim oznakama



Zalepimo čvorove za identičnim oznakama



Upareni De Brujinov graf na osnovu datog genoma

Kako konstruisati de upareni de Brujinov graf na osnovu uparenog k -gramskog sastava?



- Pretpostavili smo da je dat genom (niska *Genome*). Posmatrali smo genom kao putanju u grafu obeleženom na osnovu njegovog uparenog k -gramskog sastava
- Sada pretpostavimo da **nije dat** genom već samo upareni k -gramske sastav

Konstrukcija uparenog De Bruijinovog grafa na osnovu uparenih k -grama

TAA
GCC

ATG
CAT

GCC
TGG

CAT
GGA

TGG
ATG

GGA
GTT

AAT
CCA

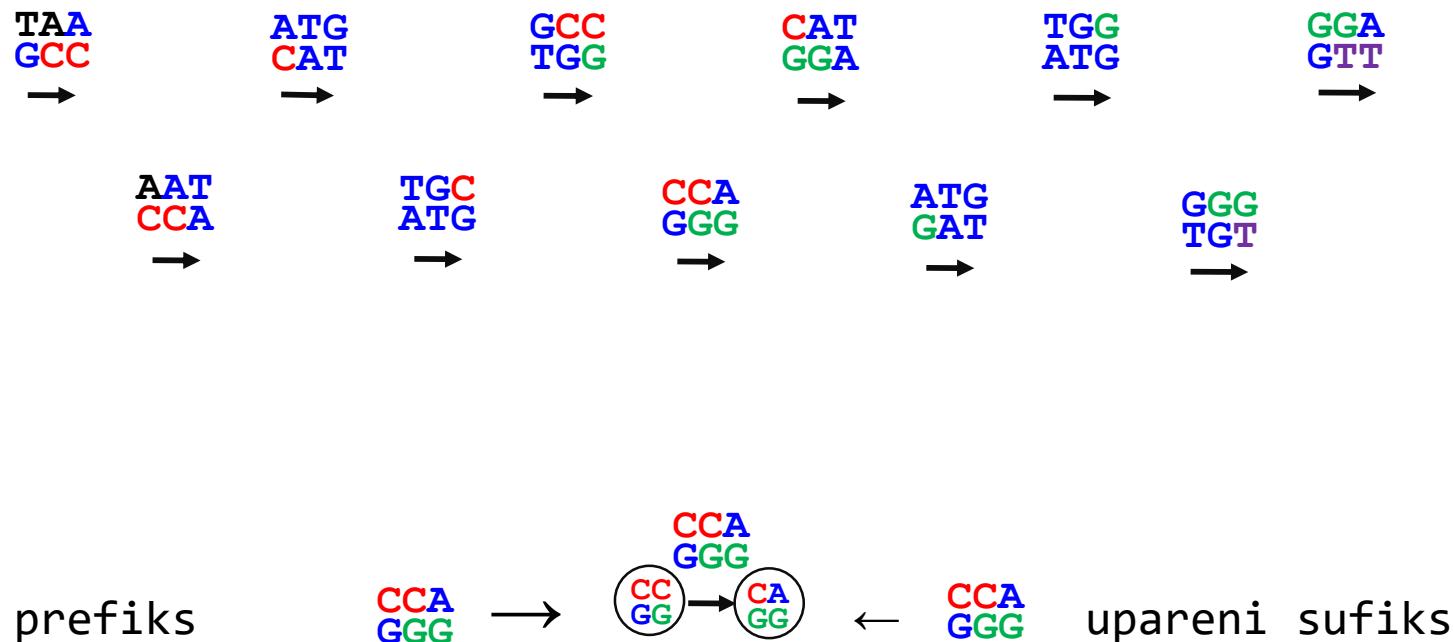
TGC
ATG

CCA
GGG

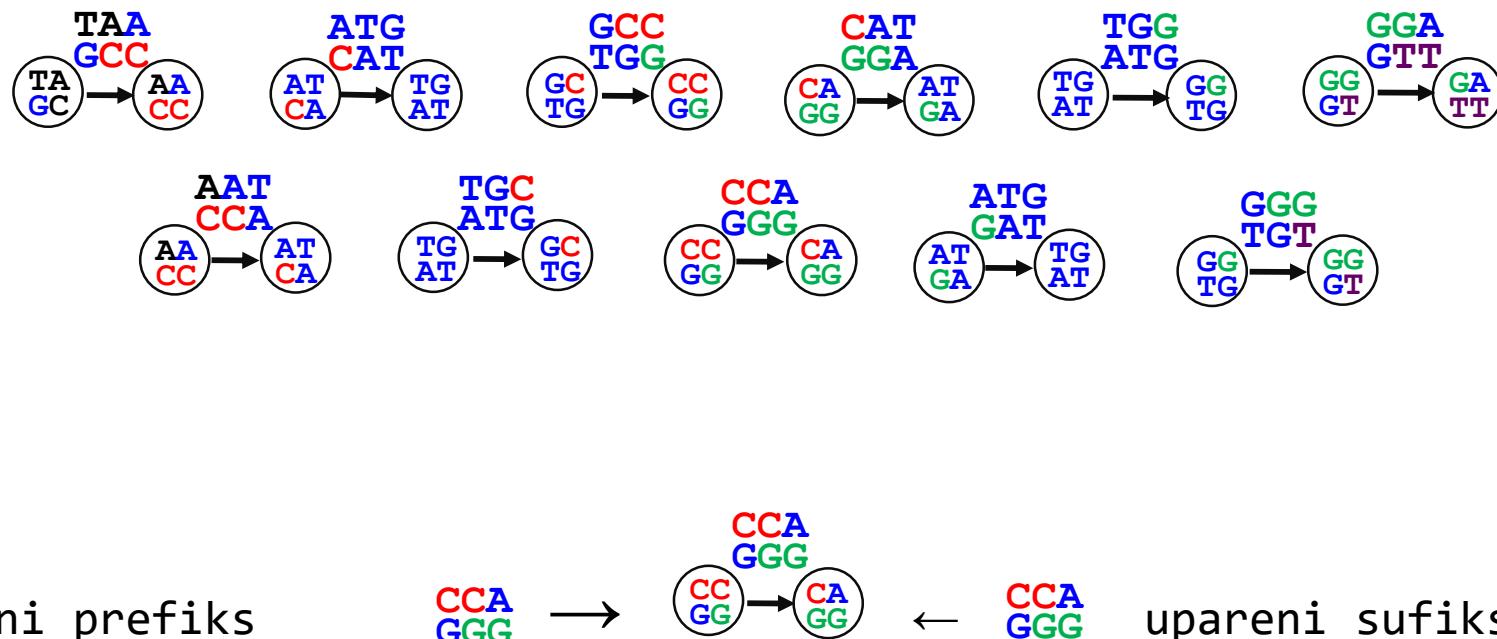
ATG
GAT

GGG
TGT

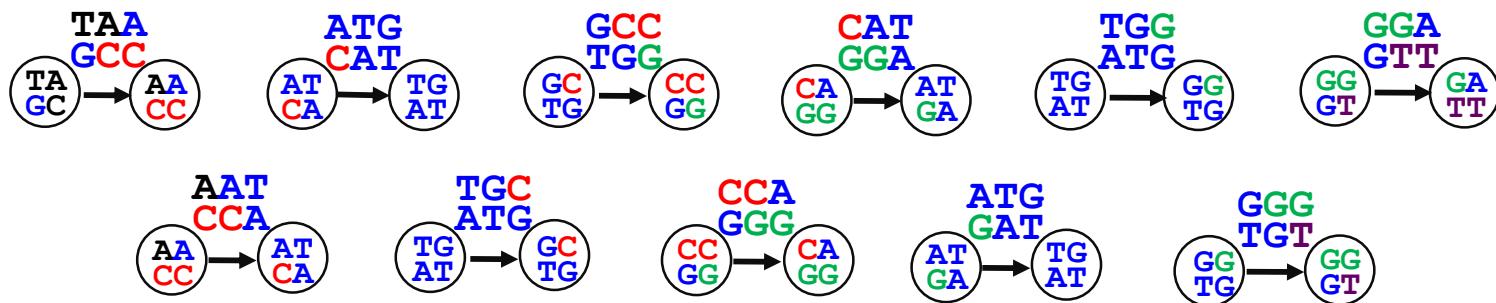
Konstrukcija uparenog De Brujinovog grafa na osnovu uparenih k -grama



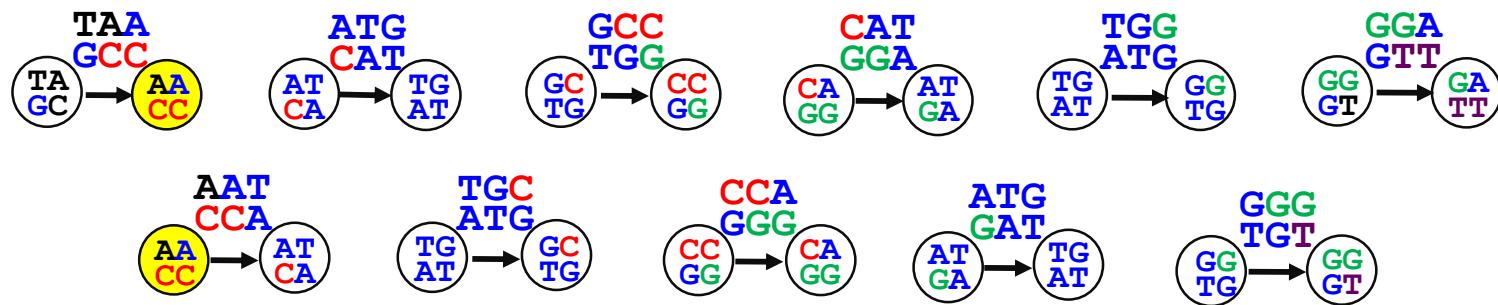
Konstrukcija uparenog De Bruijinovog grafa na osnovu uparenih k -grama



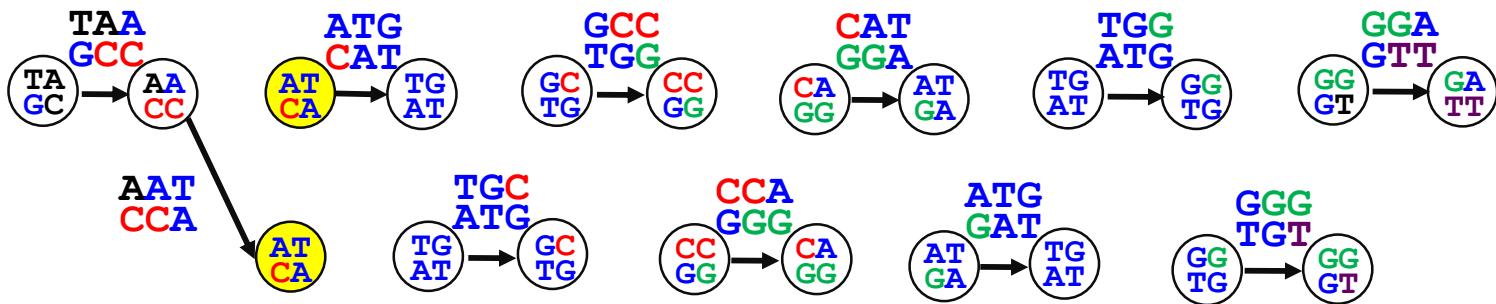
Konstrukcija uparenog De Brujinovog grafa na osnovu uparenih k -grama



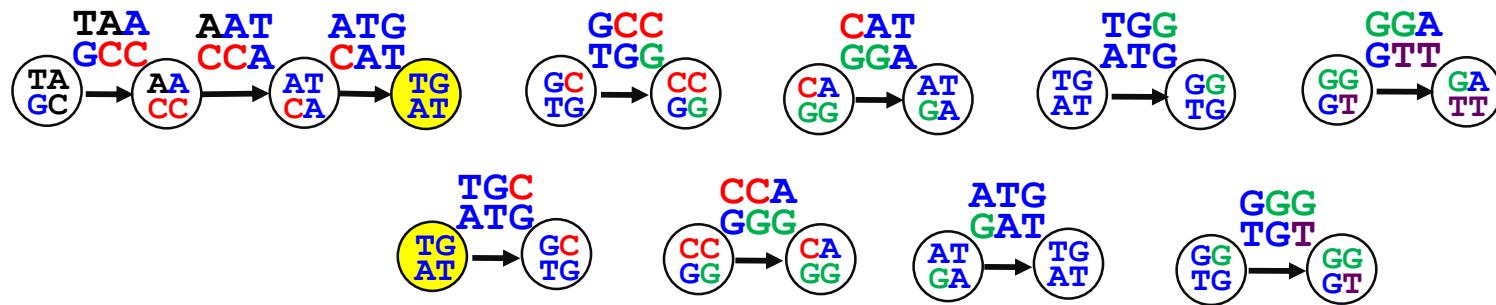
Konstrukcija uparenog De Brujinovog grafa



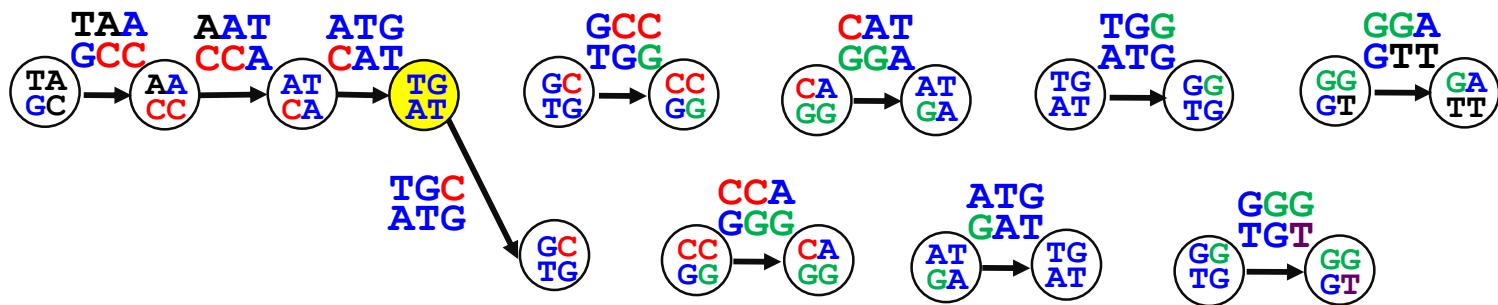
Konstrukcija uparenog De Brujinovog grafa



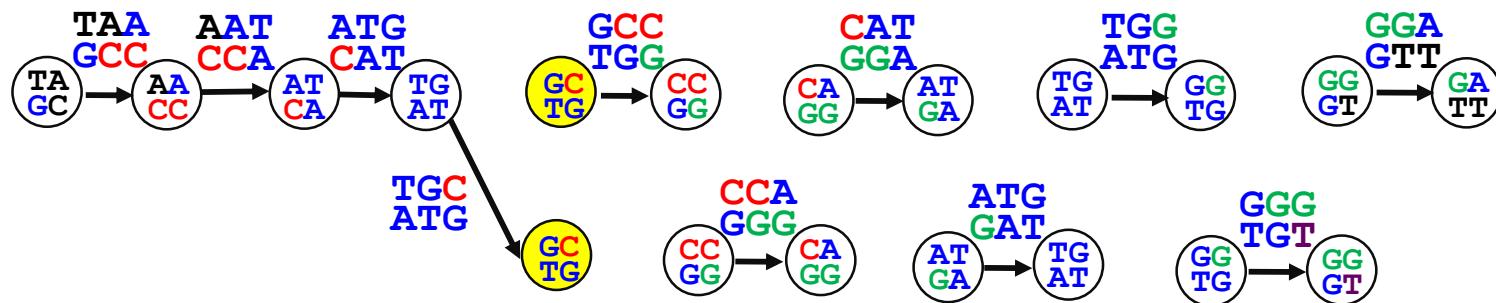
Konstrukcija uparenog De Brujinovog grafa



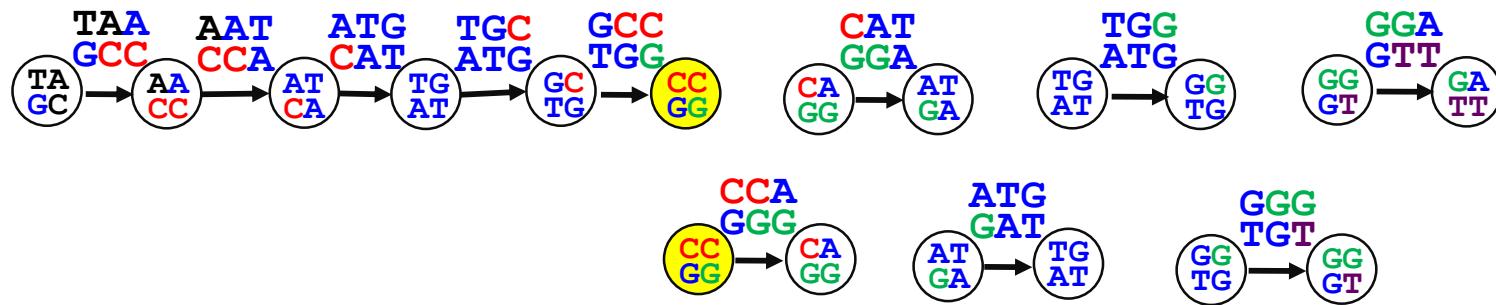
Konstrukcija uparenog De Brujinovog grafa



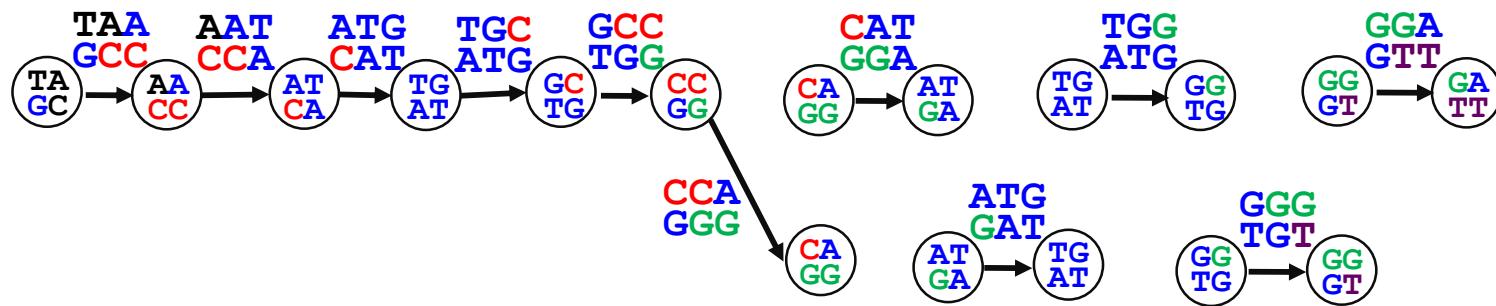
Konstrukcija uparenog De Brujinovog grafa



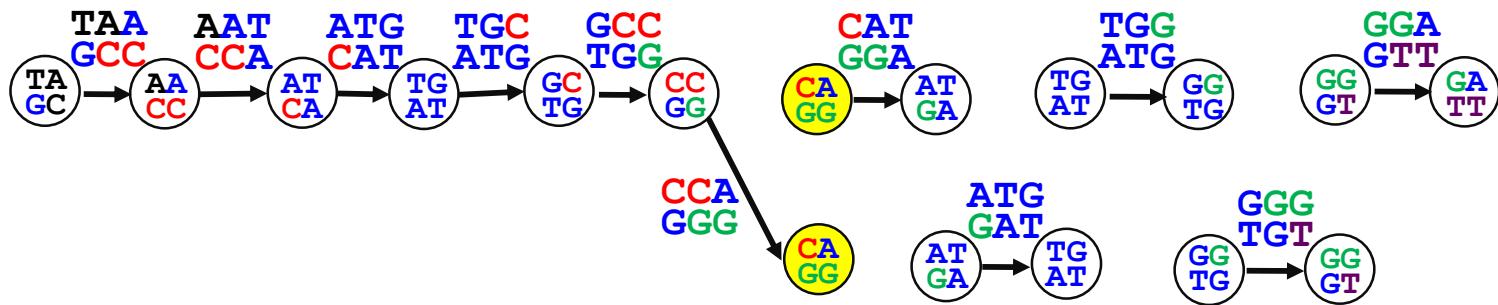
Konstrukcija uparenog De Brujinovog grafa



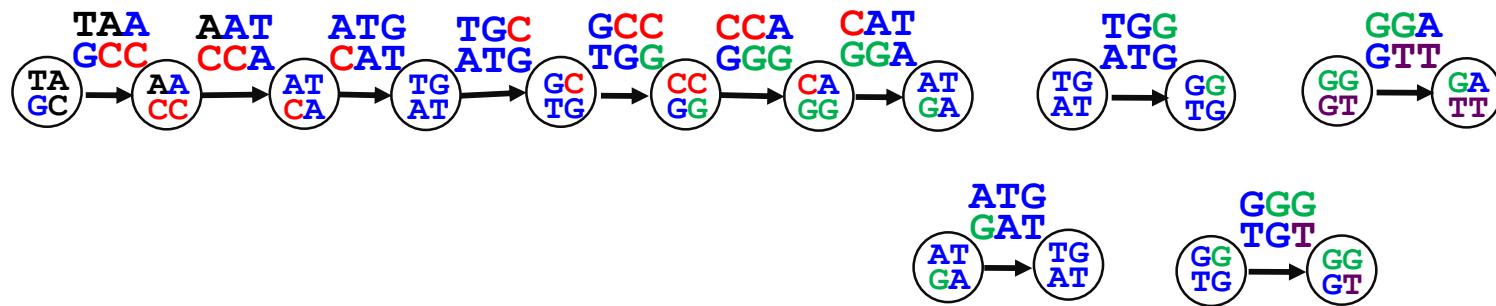
Konstrukcija uparenog De Brujinovog grafa



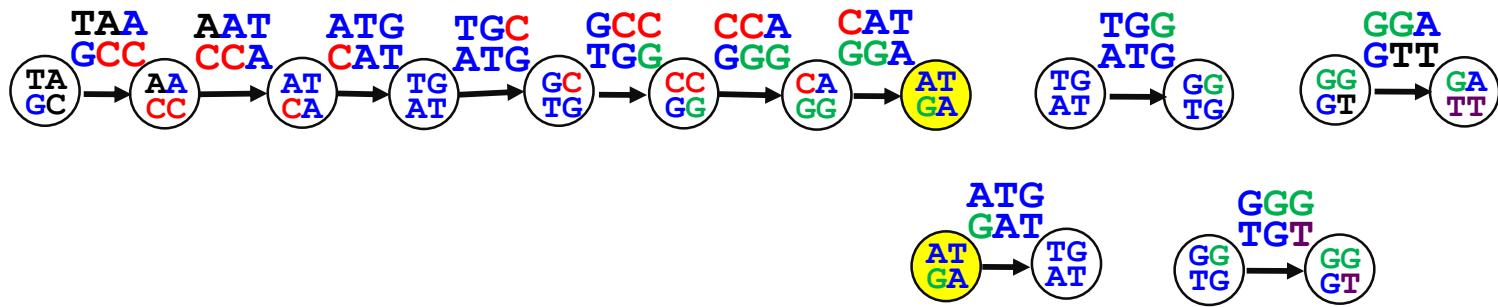
Konstrukcija uparenog De Brujinovog grafa



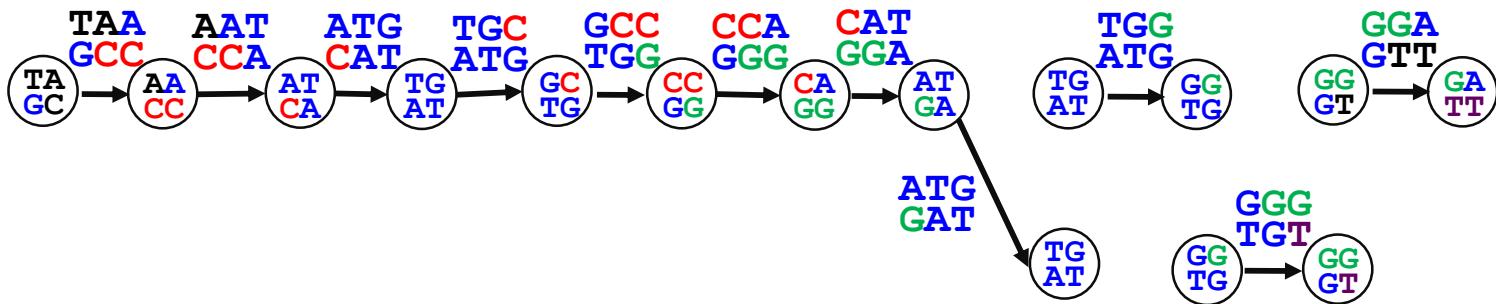
Konstrukcija uparenog De Brujinovog grafa



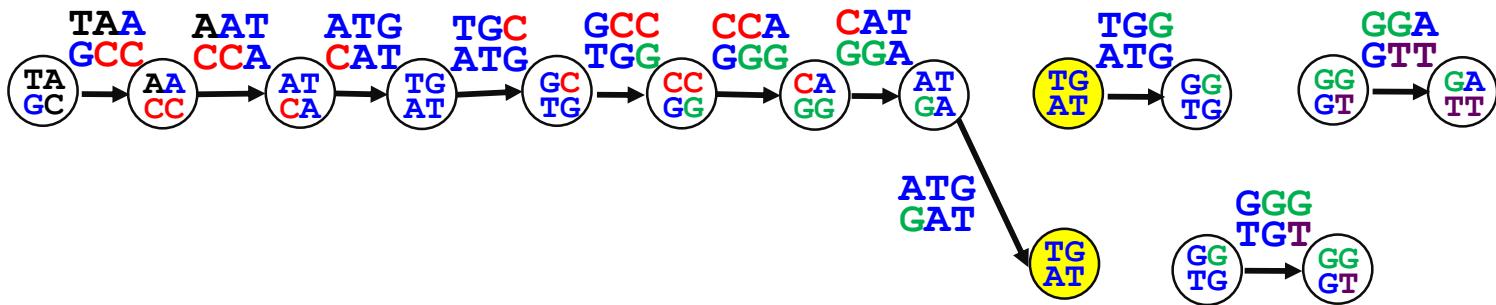
Konstrukcija uparenog De Brujinovog grafa



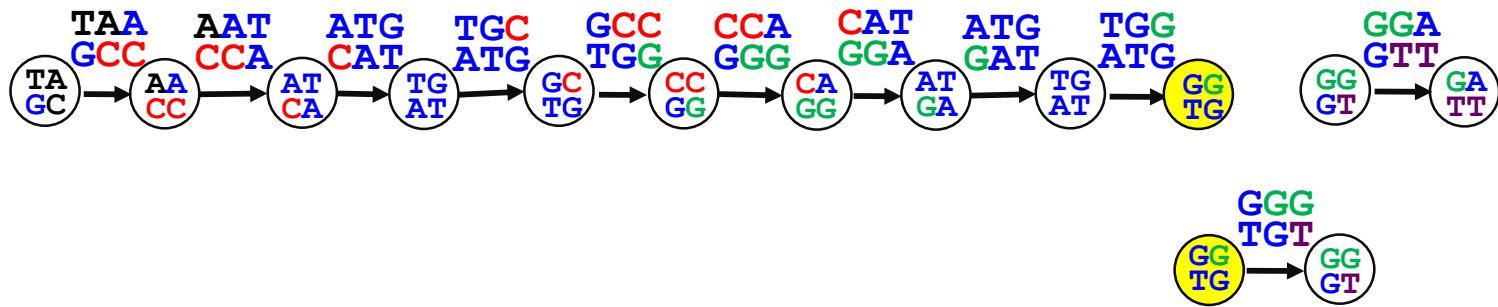
Konstrukcija uparenog De Brujinovog grafa



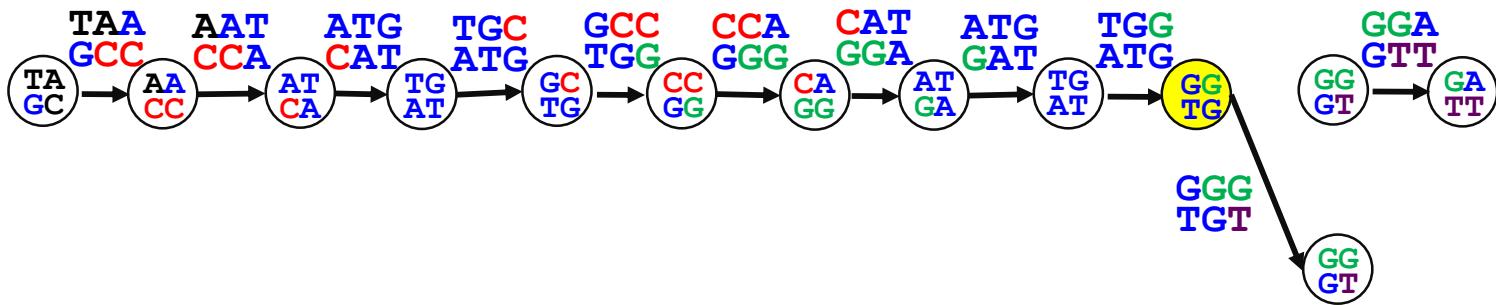
Konstrukcija uparenog De Brujinovog grafa



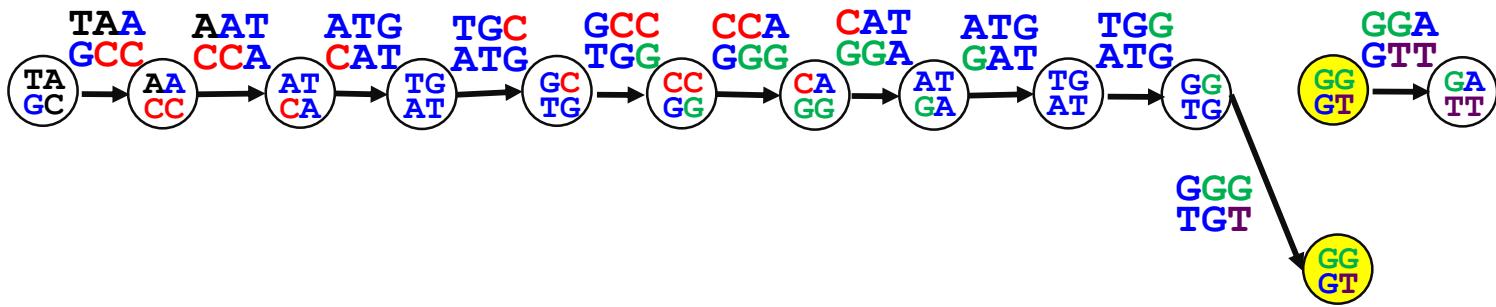
Konstrukcija uparenog De Brujinovog grafa



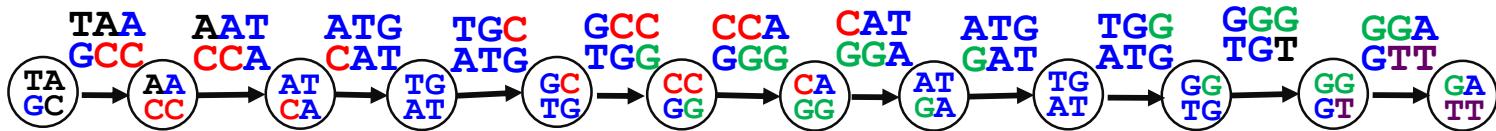
Konstrukcija uparenog De Brujinovog grafa



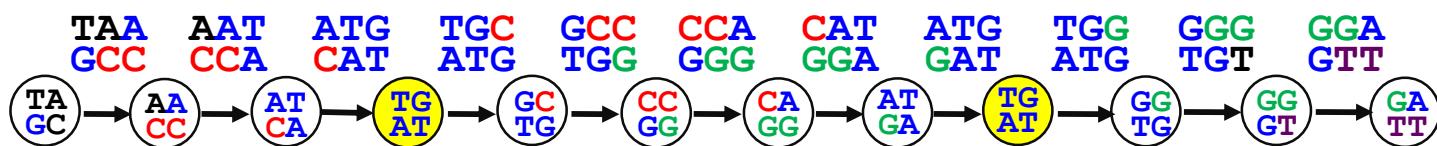
Konstrukcija uparenog De Brujinovog grafa



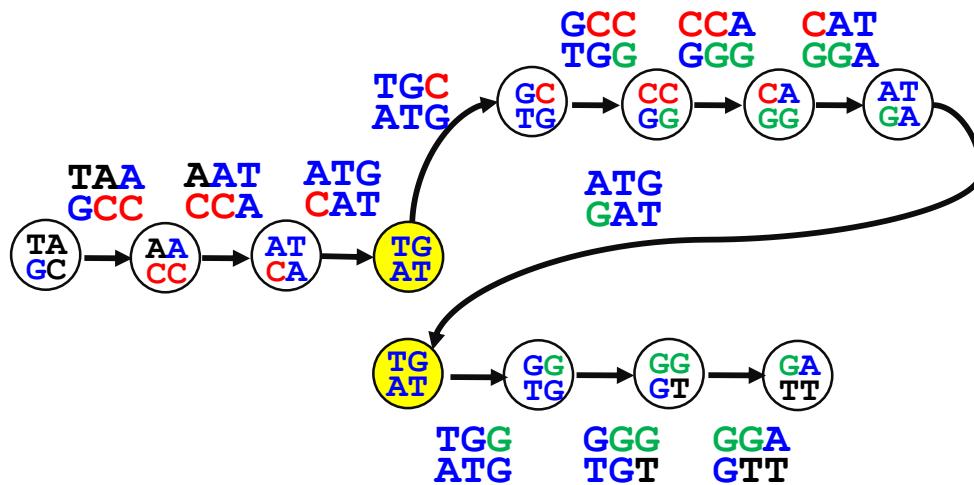
Konstrukcija uparenog De Brujinovog grafa



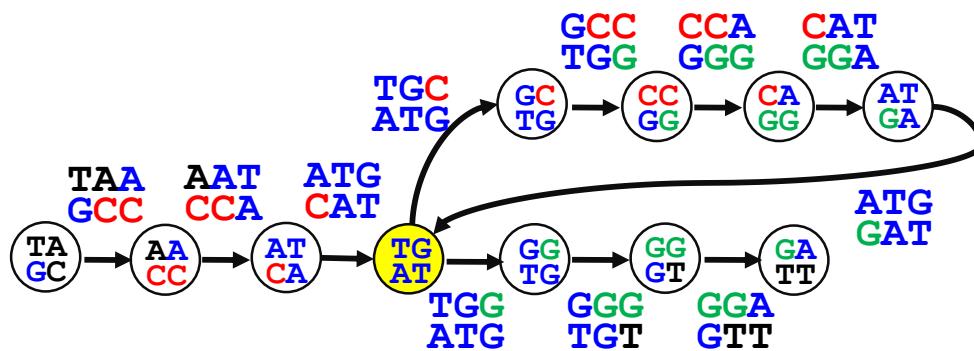
Lepljenje čvorova sa identičnom oznakom



Konstrukcija uparenog De Brujinovog grafa



Konstrukcija uparenog De Brujinovog grafa



Uparedni De Brujinov graf na osnovu parova očitavanja

Upareni De Bruijinov graf

Upareni De Bruijinov graf na osnovu kolekcije uparenih k -grama:

- Svaka grana je označena jednim uparenim k -gramom
- Svaki čvor je označen prefiksima/sufiksima izlazne/ulazne grane
- Zalepljeni su svi čvorovi sa identičnim oznakama.

Koji graf je bolja reprezentacija?

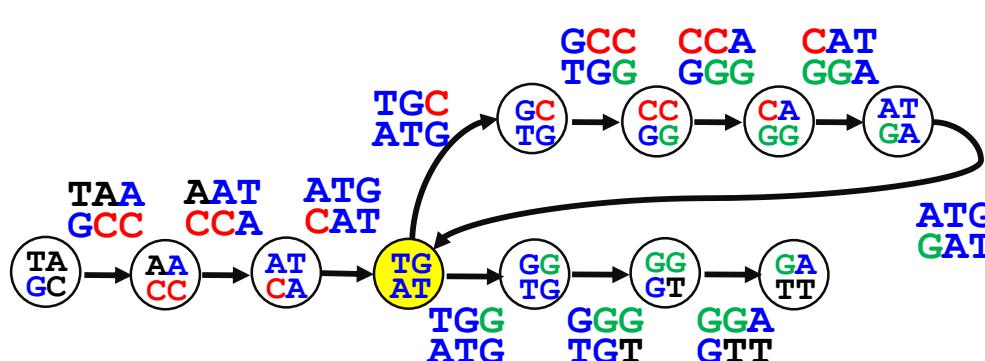
Jedinstvena
rekonstrukcija
genoma

Višestruka
rekonstrukcija
genoma

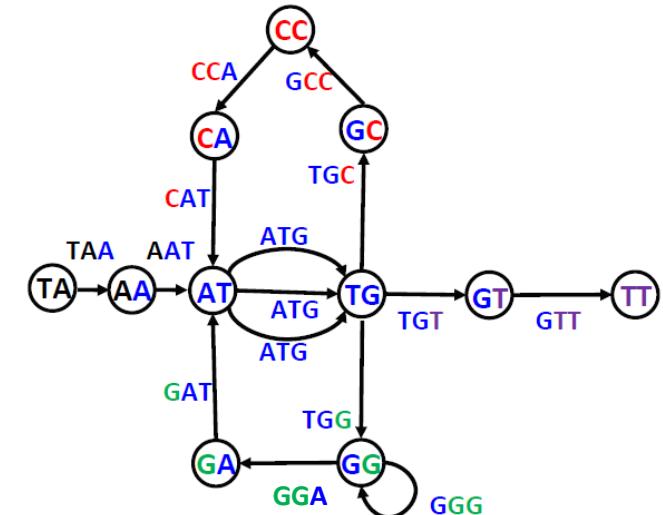
TAATG**CC**CATGGGATGTT

TAATG**CC**CATGGGATGTT

TAATGGGATG**CC**CATGTT



Upareni De Bruijinov graf



De Brujinov graf

Pregled

- Šta je sekvencioniranje genoma?
- Eksplozija u štampariji
- Problem rekonstrukcije niske
- Rekonstrukcija niske kao problem Hamiltonove putanje
- Rekonstrukcija niske kao problem Ojlerove putanje
- Slični problemi sa različitim sudbinama?
- De Bruijinovi grafovi
- Ojlerova teorema
- Sastavljanje parova očitavanja
- **U realnosti**

Nerealne pretpostavke

- Savršena pokrivenost genoma očitavanjima (svaki k-gram iz genoma je očitan)
- Očitavanja ne sadrže greške
- Rastojanja između očitavanja u okviru parova očitavanja su egzaktna

Nerealne pretpostavke

- **Nesavršena** pokrivenost genoma očitavanjima (svaki k-gram iz genoma je očitan)
- Očitavanja **ne sadrže** greške
- Rastojanja između očitavanja u okviru parova očitavanja **nisu** egzaktna
- **Itd.**

Prva nerealna pretpostavka: savršena pokrivenost

atgccgtatggacaacgact

atgccgtatg

 gccgtatgga

 gtatggacaa

 gacaacgact

Očitavanja dužine 250 nukleotida dobijena *Illumina* tehnologijom predstavljaju samo mali deo 250-grama unutar genoma.

Rešenje: razbiti dobijena očitavanja na kraće k -grame

atgccgtatggacaacgact
atgccgtatg
 gccgtatgga
 gtatggacaa
 gacaacgact

atgccgtatggacaacgact
atgcc
 tgccg
 gccgt
 ccgta
 cgtat
 gtatg
 tatgg
 atgga
 tggac
 ggaca
 gacaa
 acaac
 caacg
 aacga
 acgac
 cgact

Druga nerealna pretpostavka: očitavanja ne sadrže greške

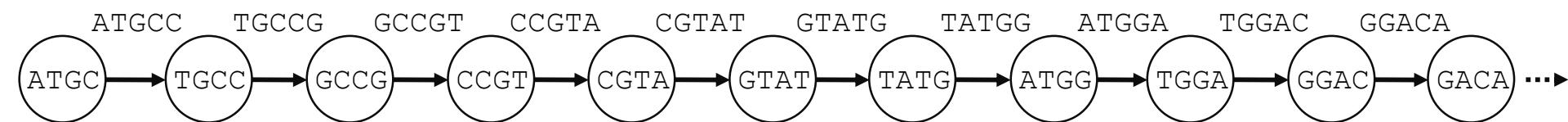
atgccgtatggacaacgact
atgccgtatg
 gccgtatgga
 gtatggacaa
 gacaacgact
 cgtacggaca

Očitavanje sa
greškom (promena
t u C)

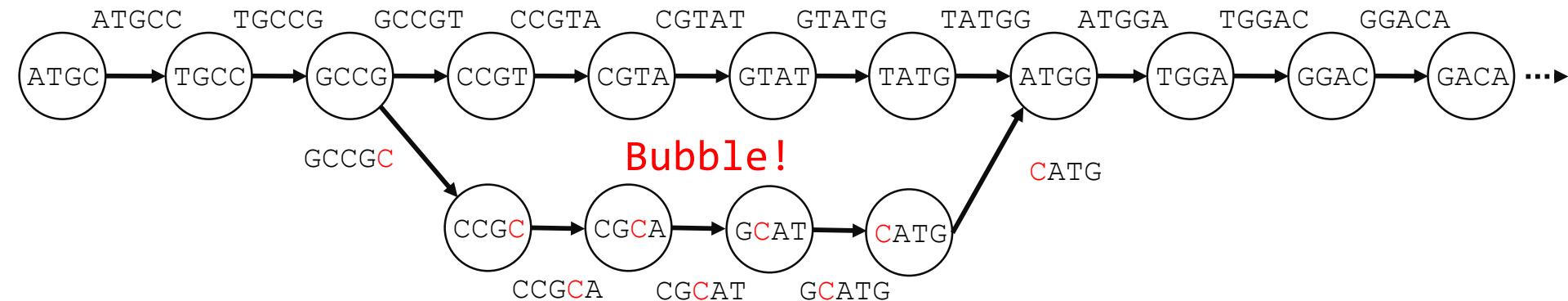
atgccgtatggacaacgact
atgcc
 tgccg
 gccgt
 ccgta
 cgtat
 gtatg
 tatgg
 atgga
 tggac
 ggaca
 gacaa
 acaac
 caacg
 aacga
 acgac
 cgact

cgtac
 gtaCg
 taCgg
 aCgga
 Cggac

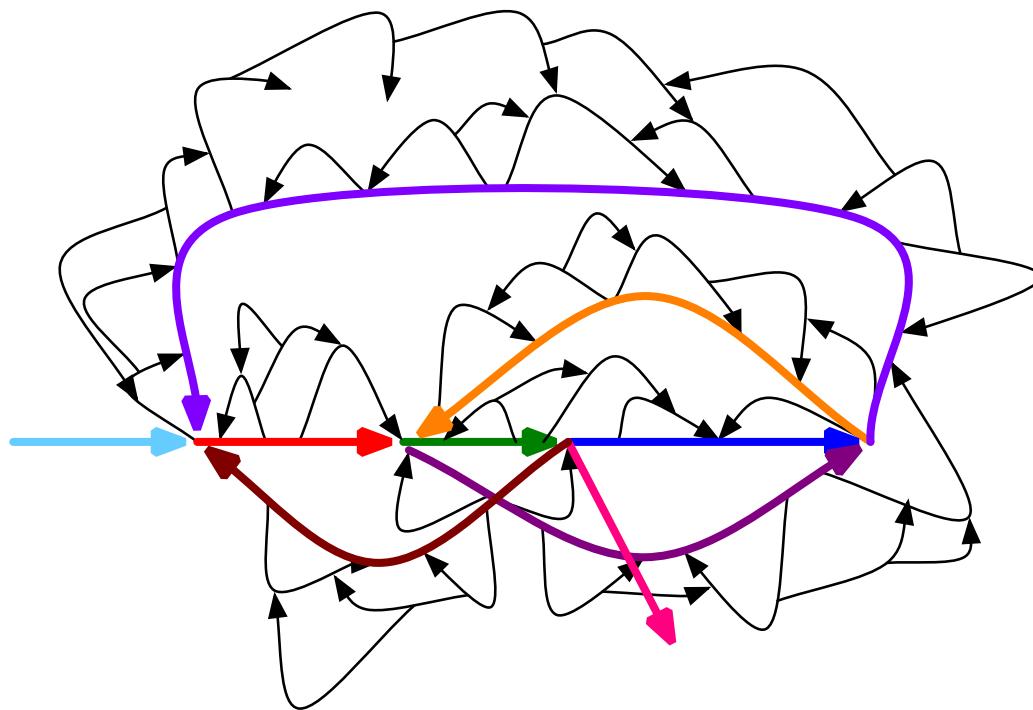
De Brujinov graf genoma ATGGCGTGCAATG... kostruisan na osnovu očitavanja koja ne sadrže greške



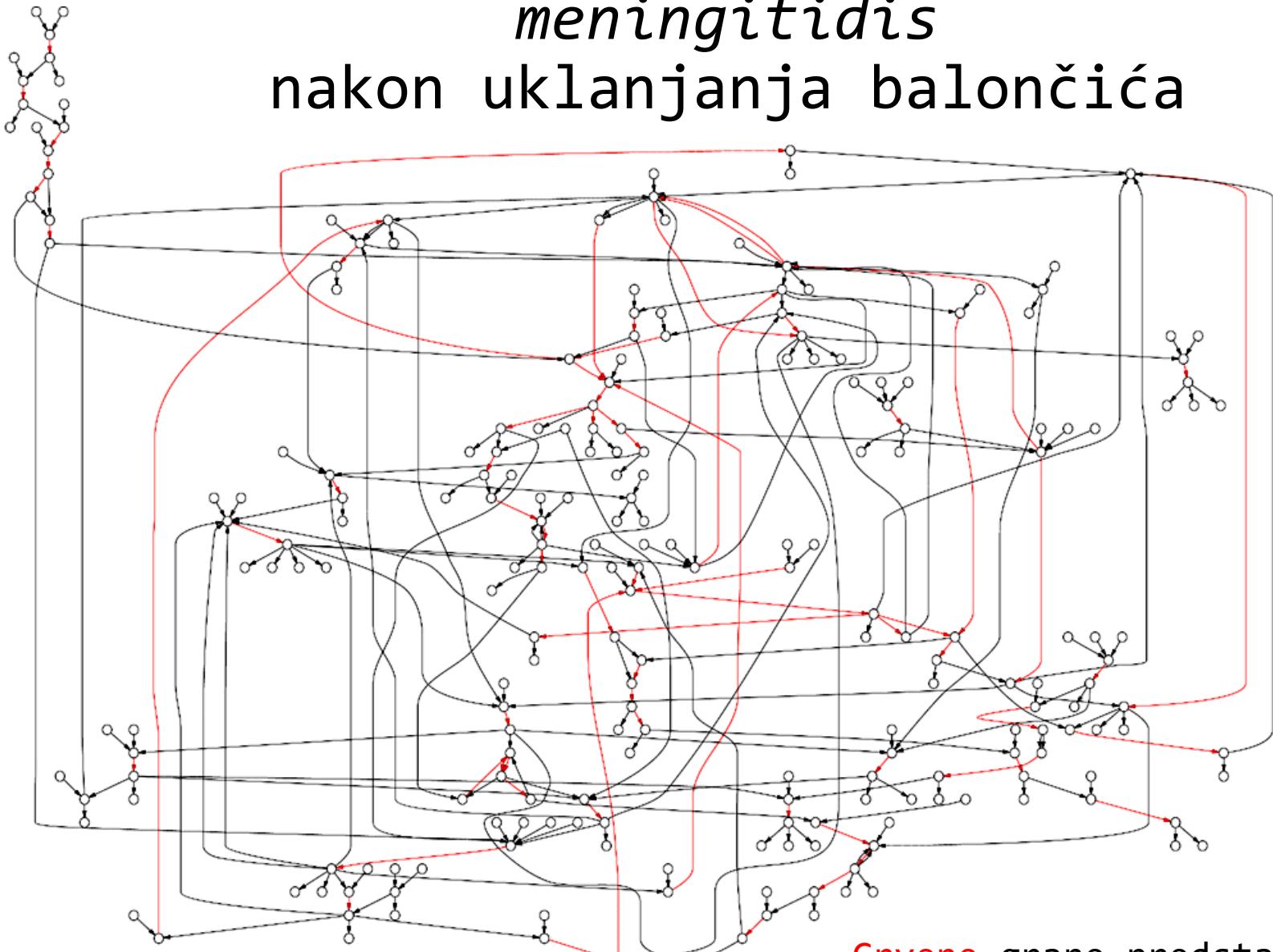
Greške u očitavanjima vode do stvaranja balončića u De Brujinovom grafu



Eksplozija balončića



De Bruijinov graf genoma *N. meningitidis* nakon uklanjanja balončića



Crvene grane predstavljaju
ponavljanja

- Slajdovi pokrivaju poglavlje 3 knjige *Bioinformatics Algorithms: an Active Learning Approach*
- Sadržaj slajdova je preuzet sa zvaničnih prezentacija autora i dodatno prilagođen