

SVEUČILIŠTE U ZAGREBU  
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

ZAVRŠNI RAD br. 6149

**De novo sastavljanje genoma  
vođeno referencom**

Sara Bakić

Zagreb, lipanj 2019.

Zagreb, 13. ožujka 2019.

## ZAVRŠNI ZADATAK br. 6149

Pristupnik: **Sara Bakić (0036498154)**  
Studij: Računarstvo  
Modul: Računarska znanost

Zadatak: **De novo sastavljanje genoma vođeno referencom**

### Opis zadatka:

Sastavljanje genoma jedan je od najsloženijih zadataka u području bioinformatike i računalne biologije. S obzirom da je problem NP težak, koriste se brojne heurističke metode. Cilj ovoga rada je testirati važnost pravilne detekcije kimernih i repetitivnih očitavanja te pronalaska svih preklapanja u de novo sastavljanju genoma. Koristeći alat minimap2 mapirati očitavanja na poznati referentni genom. Nakon toga detektirati kimerna i repetitivna očitavanja s posebnim naglaskom na repetitivnim očitavanjima koja povezuju udaljene regije genomu. Koristeći preostala očitavanja sastaviti genom pomoću alata za sastavljanje genoma Ra. Na kraju usporediti rezultate s i bez detekcije očitavanja te isprobati različite parametre za mapiranje.

Programski kod je potrebno komentirati i pri pisanju pratiti neki od standardnih stilova. Kompletnu aplikaciju postaviti na repozitorij Github.

U svezi dobivanja detaljnih informacija obratiti se Robertu Vaseru, mag. ing.

Zadatak uručen pristupniku: 15. ožujka 2019.

Rok za predaju rada: 14. lipnja 2019.

Mentor:



---

Prof. dr. sc. Mile Šikić

Djelovođa:



---

Izv. prof. dr. sc. Tomislav Hrkać

Predsjednik odbora za  
završni rad modula:



---

Doc. dr. sc. Marko Čupić

*Zahvaljujem se svima koji su savjetom i podrškom pomogli mi, kako u izradi ovoga rada, tako i u mom preddiplomskom obrazovanju uopće. Posebna zahvala mojim mentorima Mili Šikiću i Robertu Vaseru na divnoj suradnji, strpljenju i pruženoj pomoći.*

# SADRŽAJ

<b>1. Uvod</b>	<b>1</b>
<b>2. Sastavljanje genoma</b>	<b>3</b>
2.1. Ponavljajuća očitavanja . . . . .	3
2.2. Kimerna očitavanja . . . . .	4
<b>3. Podaci</b>	<b>6</b>
3.1. FASTA format . . . . .	6
3.2. FASTQ format . . . . .	7
3.3. PAF format . . . . .	8
<b>4. Metode</b>	<b>9</b>
4.1. Priprema podataka alatom minimap2 . . . . .	9
4.2. Detekcija ponavljajućih i kimernih očitavanja . . . . .	11
4.3. Potencijalni problemi prilikom detekcije kimernih i ponavljajućih očitavanja . . . . .	14
4.4. Detekcija ponavljajućih regija na referenci . . . . .	16
4.5. Izlaz programa i sastavljanje genoma . . . . .	17
<b>5. Implementacija i korištenje</b>	<b>18</b>
<b>6. Testiranje i rezultati</b>	<b>20</b>
6.1. Daljnji rad . . . . .	33
<b>7. Zaključak</b>	<b>34</b>
<b>Literatura</b>	<b>35</b>

# 1. Uvod

Bioinformatika (*bios = život + informatika*) interdisciplinarno je znanstveno područje koje povezuje biološke podatke s tehnikama pohrane, distribucije i analize informacija u svrhu istraživanja u brojnim područjima, na čelu s biomedicinom. Obuhvaća biologiju, računarsku znanost, matematiku i statistiku te ih koristi u procesu analize i interpretacije bioloških podataka. Veliki napredak u tehnikama generiranja bioloških podataka te drastičan pad u cijeni procesiranja tih podataka, omogućio je ubrzan razvoj bioinformatike (Lesk, 2013).

Obzirom da je sva genetska informacija sadržana u DNA ili RNA lancu, informacije potrebne bioinformatičari su upravo podaci o slijedu nukleotida u tim lancima. Zbog toga jedni od temeljnih oblika bioloških podataka u bioinformatičkim istraživanjima su sekvence genoma ili genomi u cijelosti.

Sekvenciranje je proces određivanja poretka nukleinskih baza – adenina (A), citozina (C), gvanina (G), timina (T), uracila (U) – u DNA ili RNA lancu. Pomoću tih informacija raznim se tehnikama utvrđuju genomi, srodnosti i evolucije vrsta, novi geni koji se pridružuju bolestima te se pokušavaju pronaći adekvatni lijekovi.

Postoji nekoliko metoda sekvenciranja koje kao zajedničko obilježje imaju ograničenu sposobnost čitanja nukleotida. Ovisno o metodi, duljine pročitanih nukleotida kreću se između pedeset i nekoliko stotina tisuća nukleotida. Kako bi sekvenciranje bilo uspješno nužno je podijeliti lance na manje fragmente te onda primjenom neke metode sekvenciranja očitati slijed nukleotida svakoga fragmenta. Dominantom strategijom danas se smatra *shotgun* sekvenciranje kojim se DNA umnaža i potom na slučajan način lomi u mnogo malih fragmenata koji se potom sekvenciraju. Obzirom da je krajnji cilj imati potpuno sastavljenu sekvencu, kada bi na ovaj način samo jednom razbili DNA na mnoštvo manjih fragmenata bilo bi nemoguće ponovno odrediti stvarni poredak tih fragmenata i sastaviti sekvencu. Upravo zato, *shotgun* sekvenciranje više se puta ponavlja kako bi se dobila očitavanja iz kojih je pomoću preklapajućih krajeva moguće sastaviti kontinuiranu sekvencu.

U okviru *shotgun* sekvenciranja postoje dva temeljna pristupa.

Prvi pristup je hijerarhijsko *shotgun* sekvenciranje koje kao prvi korak podrazumijeva izradu fizičke mape genoma. Pomoću te mape određuje se minimalan broj segmenata potreban za sekvenciranje genoma u cijelosti. Genom se potom dijeli u dugačke fragmente parova baza te se ti fragmenti kloniraju umetanjem u bakterije domaćine koristeći BAC-ove (*Bacterial artificial chromosome*). Naposljetku se dobivene kopije fragmenata *shotgun* metodom usitnjavaju na manje dijelove i sekvenciraju.

Problem kod ovakvog pristupa je izrada fizičke mape genoma jer je to skup i spor proces. Umetanje BAC-ova u bakterije i čekanje da se one razmnože kako bi se onda ti fragmenti izdvajali iz bakterija zahtijeva puno vremena te stručno osoblje.

Iz tih razloga razvijen je drugi pristup, *shotgun* sekvenciranje cijelog genoma. DNA se dijeli u slučajne dijelove raznih duljina koji se potom kloniraju. Fragmenti se zatim sekvenciraju i to uvijek od 5' kraja prema 3' te se rekonstrukcijom očitavanja dobiva potpuni slijed.

Obzirom da *shotgun* sekvenciranje cijelog genoma ima veliku vremensku i cjenovnu prednost nad hijerarhijskim *shotgun* sekvenciranjem, dominantno je u današnje vrijeme. Ipak, postoji jedan problem s takvim sekvenciranjem, a to je mogućnost ispravnog sastavljanja genoma, pogotovo onih genoma koji imaju ponavljajuće regije (Šikić i Domazet-Lošo, 2013).

Temeljni fokus u ovom radu je na sastavljanju genoma vođenom referentnom sekvencijom.

U drugom poglavlju nalazi se teoretska pozadina problema i pojašnjenja pojmova koja su fokus ovoga rada te koja će se kroz cijeli rad pojavljivati.

U trećem poglavlju prikazani su i ukratko pojašnjeni formati podataka koji se koriste u bioinformatici, a koje sam koristila u svome radu.

U četvrtom poglavlju potanko su objašnjeni algoritmi i korištene metode.

U petom poglavlju nalaze se opis i objašnjenje implementacije.

Šesto poglavlje prikazuje rezultate testiranja.

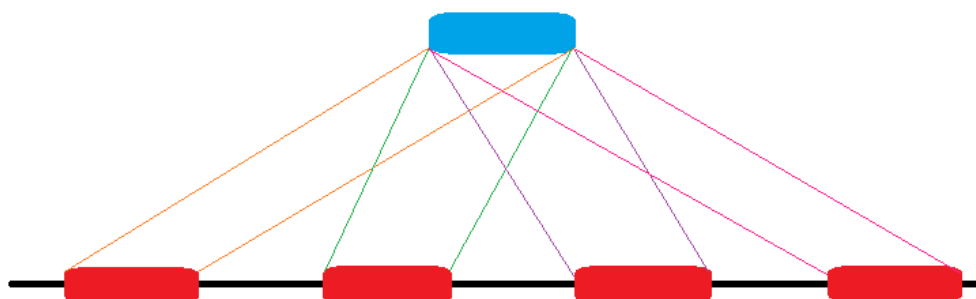
## 2. Sastavljanje genoma

Sastavljanje genoma jedan je od najsloženijih problema u području bioinformatike. Pripada skupu NP teških problema i zbog toga se u rješavanju koriste brojne heurističke metode. Problem u sastavljanju genoma predstavljaju takozvana ponavljajuća i kimerna očitavanja.

### 2.1. Ponavljajuća očitavanja

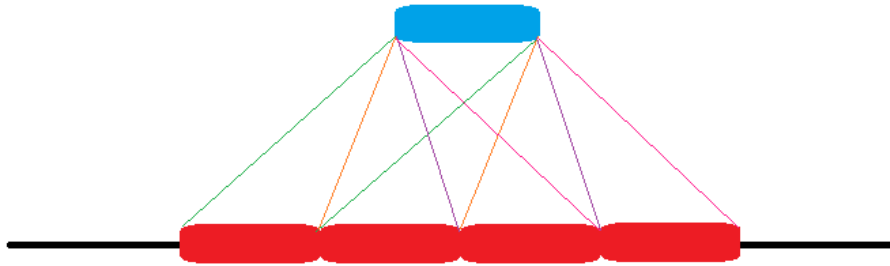
Ponavljajućim očitanjima smatramo ona očitavanja koja se mogu mapirati na više mjesta u referentnom genomu zbog postojanja dugačkih istih ili vrlo sličnih sljedova u genomu. Postoje dvije temeljne vrste ponavljajućih regija (Šikić i Domazet-Lošo, 2013).

Prva vrsta su one regije koje se nalaze na udaljenim područjima u genomu i nisu susjedne te se takve regije nazivaju raspršene (interspersed) ponavljajuće regije. Raspršene ponavljajuće sekvence najčešće nastaju kao produkt takozvanih *transposona*, odnosno dijelova DNA koji, uz pomoć enzima, imaju mogućnost umetanja na određena mjesta u kromosomu.



**Slika 2.1:** Raspršene ponavljajuće regije na referenci

S druge strane, postoje takozvane tandem ponavljajuće regije koje se u genomu pojavljuju jedna za drugom, grupirane na jednom dijelu genoma. One nastaju kao posljedica raznih bioloških mehanizama te se pojavljuju u različitim duljinama i brojnostima unutar genoma (Hauth i Joseph, 2002).



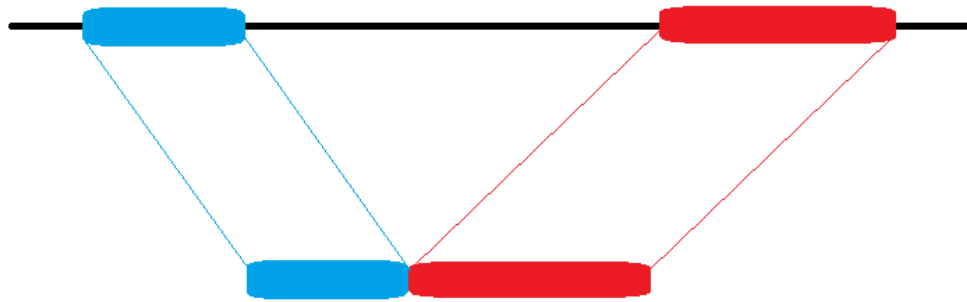
**Slika 2.2:** Tandem ponavljajuće regije na referenci

Gore objašnjen *shotgun* pristup sekvenciranju cijelog genoma u kombinaciji s ponavljajućim regijama genoma, posebice onima duljima, stvara potencijalne probleme u sastavljanju genoma. Ako prilikom *shotgun* sekvenciranja izgeneriramo sekvence koje su takve da pojedine ponavljajuće regije genoma ne mogu u cijelosti premostiti, nastaje dvosmislenost prilikom odabira sekvence koja stvarno pripada tom dijelu genoma te kvaliteta sastavljanja genoma postaje vrlo upitna. Upravo zato je vrlo bitno kvalitetno detektirati ponavljajuća očitavanja i ponavljajuće regije na genomu.

## 2.2. Kimerna očitavanja

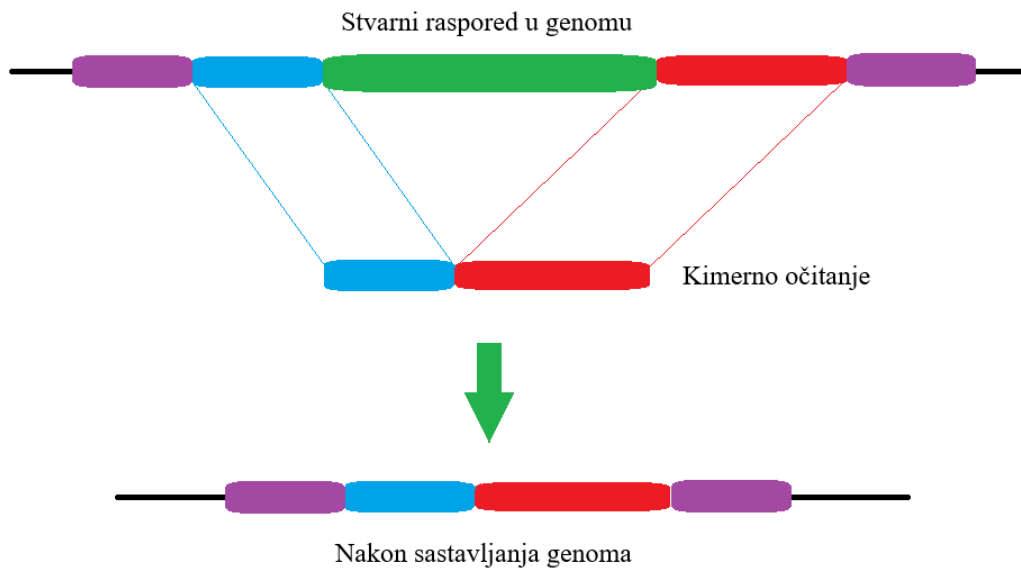
Kimernim očitavanjima (*kimera – jedinka sastavljena od dijelova različitih rasa*) smatramo ona očitavanja koja nastaju pogrešnim spajanjem dva očitavanja koja predstavljaju u prirodi veoma udaljene dijelove genoma. Oni nastaju kao posljedica grešaka prilikom identificiranja nukleotida u očitavanjima. Zbog brzine kojom uređaji novijih generacija to rade događaju se greške u očitavanjima zbog kojih se pojavljuju lažna preklapanja između prirodno vrlo udaljenih očitavanja te se kao posljedica ta očitavanja spajaju u jedinstveno očitavanje koje bi trebalo predstavljati jedan dio genoma, ali zapravo predstavlja dva vrlo udaljena dijela genoma (atdbio, 2011).





**Slika 2.3:** Kimerno očitane nastala spajanjem dvaju očitane

I takva očitane stvaraju probleme pri izgradnji genoma. Obzirom da je proces sastavljanja genoma zapravo proces traženja preklapanja krajeva očitane, kimerno očitane će na jednom kraju pronaći preklapanje s jednim dijelom genoma, a na drugoj s potpuno drugim dijelom genoma što očito narušava sliku stvarnog genoma. (Institute, 2003)



**Slika 2.4:** Problemi sa sastavljanjem genoma uzrokovani kimernim očitanjima

Osnovni cilj moga rada jest detekcija kimernih i ponavljajućih očitane kako bi, uz izostavljanje tih očitane, što kvalitetnije i točnije sastavili genom.

## 3. Podaci

### 3.1. FASTA format

Bioinformatički podaci najčešće se nalaze u FASTA tekstualnom formatu. FASTA je tekstualni format kojim se prikazuju sljedovi nukleotida, a svaki nukleotid predstavlja jedno slovo (Šikić i Domazet-Lošo, 2013).

Slovo reprezentat	Nukleotid
A	Adenin
C	Citozin
G	Gvanin
T	Timin
U	Uracil
R	Adenin ili gvanin
Y	Citozin, timin ili gvanin
K	Gvanin, timin ili uracil
M	Adenin ili gvanin
S	Citozin ili gvanin
W	Adenin, timin ili uracil
B	Ne adenin
D	Ne citozin
H	Ne gvanin
V	Niti timin niti uracil
N	Adenin, citozin, gvanin, timin, uracil
X	Maskiranje
-	Procjep neodređene duljine

**Tablica 3.1:** Tablica znakova u FASTA formatu

Svaka sekvenca u FASTA formatu sastoji se od dvije linije. Prva linija započinje znakom „>“ i nakon nje slijedi identifikator slijeda. Druga linija sastoji se od samoga slijeda.

```
>a019fb87-85b7-495a-b6dc-789a2f8c4572_Basecall_2D_000_2d  
TACGCATAAGCGCCAAAAGCACAAGATGCTCACCGCCAG...
```

## 3.2. FASTQ format

Drugi česti format je FASTQ format koji je vrlo sličan FASTA formatu.

Zapis sekvence sastoji se od 4 linije. Prva linija započinje znakom „@“ nakon koje slijedi identifikator slijeda. U drugoj liniji nalazi se sami slijed. Treća linija započinje znakom „+“ nakon kojeg opcionalno slijedi identifikator slijeda i u četvrtoj liniji nalazi se vrijednost kvalitete slijeda iz druge linije (Šikić i Domazet-Lošo, 2013).

```
@cluster_2:UMI_ATTCCG  
TTTCCGGGGCACATAATCTTCAGCCGGGCGC...  
+  
9C;=;<9@4868>9:67AA<9>65<=>591
```

Početni podaci su se nalazili upravo u FASTA i FASTQ formatu te su kao takvi predstavljali ulaz alatu minimap2 koji je na temelju očitavanja koja se najčešće nalaze u FASTQ formatu i reference koja je najčešće u FASTA datoteci mapirao očitavanja na sekvencu i to prikazao u PAF formatu.

### 3.3. PAF format

PAF format tekstualni je format koji opisuje aproksimativne pozicije mapiranja između dva seta sekvenci. Sastoji se od ovih informacija odijeljenih tabulatorom (Li, 2018b):

Stupac	Tip	Opis
1	string	Naziv očitavanja
2	int	Duljina očitavanja
3	int	Početna pozicija mapiranja očitavanja
4	int	Krajnja pozicija mapiranja očitavanja
5	char	Relativni slijed "+" ili "-"
6	string	Naziv reference
7	int	Duljina reference
8	int	Početna pozicija mapiranja očitavanja na referenci
9	int	Krajnja pozicija mapiranja očitavanja na referenci
10	int	Broj poklapanja očitavanja i reference
11	int	Ukupan broj poklapanja, promašaja, umetanja i brisanja u poravnanju
12	int	Kvaliteta mapiranja

Tablica 3.2: Opis PAF formata

```

cc8c73c3-58cb-4d3b-9ec5-acf513349500_Basecall_2D_000_2d 6325 31 6297 + NC_000913.3 4641652 2788185 2794500 3221 6494 60 tp:A:P cn:l:403 s1:l:3182 s2:l:0 dv:f
:0.0690 rl:l:27
e0ba75c5-e69a-4813-a83c-da7fed66d499_Basecall_2D_000_2d 9088 21 8865 - NC_000913.3 4641652 2201970 2210879 2849 9068 60 tp:A:P cn:l:355 s1:l:2798 s2:l:0 dv:f
:0.1031 rl:l:0
a67e5d91-7623-44db-b9bf-30d4829e9033_Basecall_2D_000_2d 9805 97 8957 + NC_000913.3 4641652 1646485 1655290 4877 9068 60 tp:A:P cn:l:665 s1:l:4818 s2:l:214 d
v:f:0.0607 rl:l:0
a1a4b08d-e16e-4a09-8a52-f6e850015994_Basecall_2D_000_2d 13498 49 13469 - NC_000913.3 4641652 120916 134061 6888 13622 60 tp:A:P cn:l:849 s1:l:6802 s2:l:0 dv:f
:0.0716 rl:l:0
21e647c8-f64a-4cbd-8402-cb9ca7efb65b_Basecall_2D_000_2d 18445 29 10430 - NC_000913.3 4641652 1014434 1024570 5293 10596 60 tp:A:P cn:l:729 s1:l:5199 s2:l:0 dv:f
:0.0655 rl:l:0
e76554aa-dbed-4317-8318-43cda9886691_Basecall_2D_000_2d 2279 58 2255 + NC_000913.3 4641652 1244442 1246644 1083 2240 60 tp:A:P cn:l:147 s1:l:1077 s2:l:0 dv:f
:0.0703 rl:l:0
90600d13-1403-491e-bd85-5c9d3aed4acb_Basecall_2D_000_2d 10703 36 10670 - NC_000913.3 4641652 1584213 1594463 5689 10729 60 tp:A:P cn:l:730 s1:l:5608 s2:l:0 dv:f
:0.0678 rl:l:0
b0f70921-2539-4bde-a9fa-45ca3e7bc9c7_Basecall_2D_000_2d 5888 38 5850 + NC_000913.3 4641652 1457128 1467841 3581 5883 60 tp:A:P cn:l:483 s1:l:3558 s2:l:0 dv:f
:0.0546 rl:l:0

```

Slika 3.1: Isječak mapiranja u PAF formatu

## 4. Metode

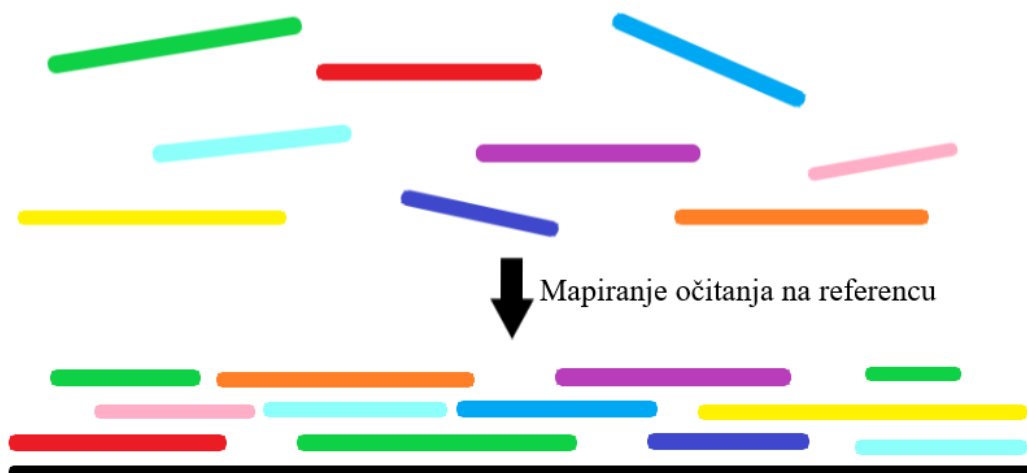
### 4.1. Priprema podataka alatom minimap2

Minimap2 svestrani je alat korišten za mapiranje DNA ili RNA sekvenci na referentnu bazu podataka.

Nekoliko je najčešćih tipova korištenja alata minimap2:

- mapiranje PacBio ili Nanopore očitavanja na ljudski genom
- traženje preklapanja između dugačkih očitavanja
- splice-aware poravnanje na referentni genom
- poravnanje jednostrukih ili uparenih Illumina očitavanja
- assembly-to-assembly poravnanje
- poravnanje cijelih genoma srodnih vrsta

Minimap2 sam koristila za mapiranje očitavanja u FASTA formatu na referentnu bazu podataka u FASTQ formatu. Kao rezultat dobila sam indeks mapiranja u PAF formatu i pomoću njega detektirala ponavljajuća i kimerna očitavanja.



Slika 4.1: Mapiranje očitavanja na referencu

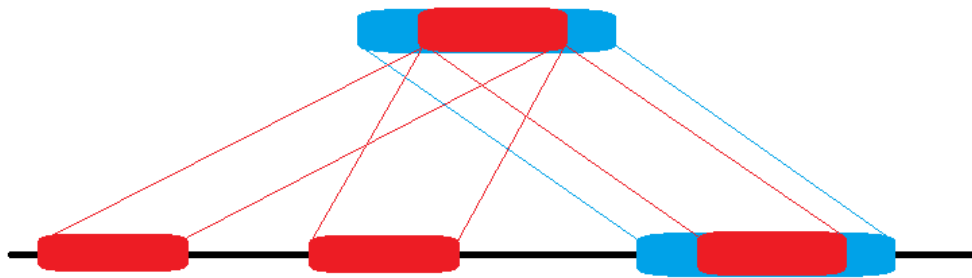
Kako bi korištenje rezultata minimap2 u svrhu detekcije kimernih i ponavljajućih očitavanja bilo što razumljivije, ovo je kratak pregled algoritma minimap2: (Li, 2018a)

1. za **I** referentnih baza se pronadū takozvani minimizeri - specifični podnizovi preddefinirane duljine **k** koji služe za detekciju sličnih regija među dvjema sekvencama - i spremaju u indeks u obliku tablice raspršenog adresiranja
2. čita se **K** baza sekvenci i za svaku se ponavljaju koraci 3 do 7
3. svaki minimizer očitavanja provjerava se u indeksu minimizera referentne baze i ako se ne nalazi među preddefiniranih **f** najčešćih, skupljaju se njegova pojavljivanja u referenci te se ona nazivaju *seeds*
4. *seeds* sortiraju se po poziciji na referenci te se dinamički povezuju u lance
5. za svaki lanac, počevši od najboljeg po njegovoj lančanoj vrijednosti, provjerava se preklapanje s lancima u postojećem (početno praznom) setu primarnih mapiranja; ako je preklapanje manje od preddefinirane vrijednosti **mask-level** dodaje se u set primarnih mapiranja, inače se deklarira kao sekundarno mapiranje lancu s kojim ima preklapanje veće od zadane vrijednosti
6. uzimaju se sva primarna mapiranja i **N** najboljih sekundarnih očitavanja koja imaju lančanu vrijednost veću od **p** posto odgovarajućeg primarnog mapiranja
7. filtriraju se oni *seedovi* koji vode do velikih umetanja i brisanja prilikom poravnanja te se nad preostalima vrši globalno poravnanje; lanci se dijele na manje ako vrijednost poravnanje padne za **z**, zanemarujući dugačke praznine te se lanci i njihova mapiranja ispisuju u odgovarajućem formatu (u našem slučaju PAF formatu)
8. ako postoji više sekvenci vraća se na korak 2 dok god postoji još neobrađenih sekvenci
9. ako postoji više referentnih sekvenci ponavlja se cijeli postupak dok god sve referentne sekvence nisu obrađene

## 4.2. Detekcija ponavljajućih i kimernih očitavanja

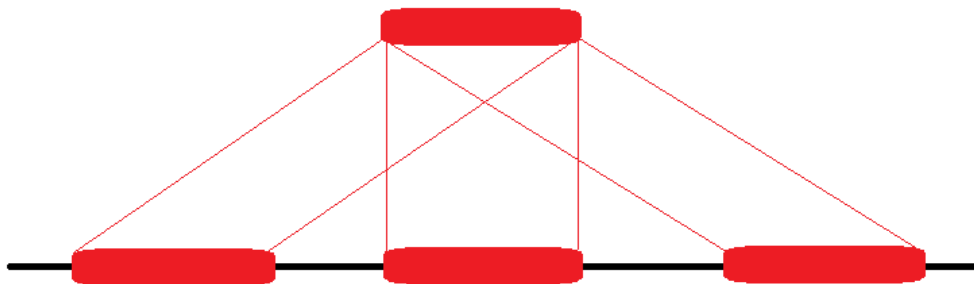
Obzirom da su ponavljajuća očitavanja definirana kao ona očitavanja koja se mogu mapirati na više mjesta na genomu, algoritam minimap2 prilikom izgradnje indeksa mapiranja detektira sva višestruka mapiranja ponavljajućih sekvenci.

Sekvenca može biti ponavljajuća u dva smisla. Primjerice, sekvenca kao takva može imati mapiranje na referentni genom cijelom svojom duljinom (misli se na barem 90% stvarne duljine sekvence) te unutar sebe sadržavati regiju koja se naknadno u jednom ili više navrata mapira na druga mjesta na genomu.



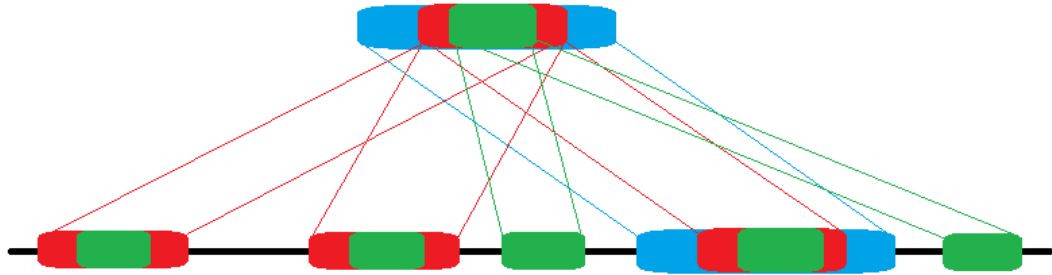
**Slika 4.2:** Ponavljajuća sekvenca s jednim primarnim mapiranjem na referencu

S druge strane, postoje sekvence koje su u cijelosti ponavljajuće, odnosno njihova najbolja (najdulja) mapiranja su višestruka.



**Slika 4.3:** Sekvenca u cijelosti ponavljajuća mapirana na referencu

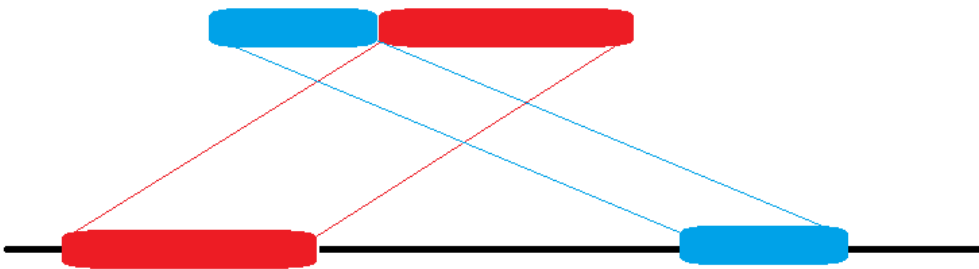
Potencijalno, unutar svake ponavljajuće regije na sekvenci može postojati manja regija koja je također ponavljajuća.



**Slika 4.4:** Ponavljajuća regija očitavanja unutar ponavljajuće regije

U gore pojašnjenom algoritmu minimap2 za ponavljajuće regije posebno su važni koraci 5 i 6 gdje se ponavljajuće regije prepoznaju kao sekundarna mapiranja odgovarajućem primarnom mapiranju te se u 6. koraku najveće ponavljajuće regije uzimaju u obzir prilikom poravnanja sekvence na referencu.

Dok ponavljajuća očitavanja možemo zamisliti kao učahurene ponavljajuće regije unutar cijelog očitavanja, temeljna karakteristika kimernih očitavanja jest da su ona "isjeckana" i tako odvojeno mapirana na različite dijelove genoma.



**Slika 4.5:** Mapiranje kimernog očitavanja na genom

Obzirom da su ona nastala pogrešnim spajanjem prirodno odvojenih očitavanja, nemoguće ih je u jednom dijelu kvalitetno mapirati na genom.

Za njih je najvažniji korak 7 u algoritmu minimap2 koji lance dijeli u manje ako vrijednost poravnanja pada, zanemarujući velike praznine.



Nakon parsiranja PAF ispisa minimap2-ovog algoritma, sva dostupna mapiranja spremljena su u mapu koja je kao ključ imala naziv sekvence koja se mapira, a kao vrijednost pridruženu ključu sva pronađena mapiranja tog očitavanja.

Kao potencijalna kimerna i ponavljajuća očitavanja nameću se ona očitavanja za koja je pronađeno više od jednog mapiranja.

Poznavajući početne i krajnje pozicije mapiranja očitavanja moguće je jasno odvojiti kimerna od ponavljajućih. Ponavljajuća očitavanja imat će višestruka mapiranja kod kojih su područja očitavanja koja se mapiraju u odnosu podskupa, dok će kimerna očitavanja imati višestruka mapiranja kod kojih su pozicije mapiranja na očitavanju praktički bez presjeka.

Primjerice, imamo očitavanje koje se mapira na šest različitih mjesta u genomu i to na način prikazan u tablici 4.1.

Početna pozicija na očitavanju	Krajnja pozicija na očitavanju	Početna pozicija na referenci	Krajnja pozicija na referenci
85	16679	1865475	1881092
7822	8577	1978527	1979239
7822	8577	257932	258644
7822	8577	19820	20532
7822	8577	3583458	3584170
7822	8577	1049808	1050520

**Tablica 4.1:** Primjer mapiranja ponavljajuće sekvence

Po pozicijama mapiranja očitavanja vidljivo je da je prvo mapiranje primarno mapiranje toga očitavanja. Gledajući ostala mapiranja lako se uočava da se sva ta mapiranja mogu zamisliti kao podregije regije koja čini primarno mapiranje zbog čega takvo očitavanje klasificiramo kao ponavljajuće očitavanje.

S druge strane moguće je i očitavanje s dva mapiranja prikazana u tablici 4.2.

Početna pozicija na očitavanju	Krajnja pozicija na očitavanju	Početna pozicija na referenci	Krajnja pozicija na referenci
17	3842	756995	760771
3884	15141	3320818	3331885

**Tablica 4.2:** Primjer mapiranja ponavljajuće sekvence

Gledajući pozicije mapiranja toga očitavanja, lako je uočljivo da pozicije prvog i drugog mapiranja nemaju presjek, odnosno prvi dio očitavanja mapira se na jedan dio reference, a drugi, odvojeni dio očitavanja mapira se na drugi dio reference. Zbog toga, takvo očitavanje klasificiramo kao kimerno očitavanje.

### 4.3. Potencijalni problemi prilikom detekcije kimernih i ponavljajućih očitavanja

Nekoliko je slučajeva u kojima nije odmah vidljivo je li neko očitavanje kimerno, ponavljajuće ili regularno očitavanje.

Prvi takav slučaj je očitavanje koje ima ovakva mapiranja:

Početna pozicija na očitanju	Krajnja pozicija na očitanju	Početna pozicija na referenci	Krajnja pozicija na referenci
170	2445	1870705	1872977
2531	3306	3583433	3584182
2531	3306	257920	258669
2531	3306	19808	20557
2531	3306	1978515	1979264
2531	3306	1049783	1050532
2531	3306	279175	279924

**Tablica 4.3:** Primjer mapiranja kimernog očitavanja s ponavljajućom regijom

Ovo očitavanje ima dvije vrste mapiranja. Vidljivo je da postoje mapiranja zbog kojih bi očitavanje detektirali kao kimerno (dijelovi očitavanja mapiraju se na pozicijama 170-2445, 2531-3306), ali postoji i regija koja ima ponavljajuće mapiranje (na pozicijama 2531-3306). Takvo mapiranje se usprkos ponavljajućoj regiji detektira kao kimerno očitavanje jer je iz mapiranja očito da je očitavanje neregularno u smislu nepripadnosti jednom dijelu genoma. Ponavljajuća regija samostalno je očitavanje koje je zapravo ponavljajuće, ali je zbog greške u sekvenciranju spojeno s drugim očitanjem u kimerno.

Drugi problem vezan je uz prokariotske kromosome koji su tipično cirkularni. Zbog cirkularnosti kromosoma, postoje očitavanja koja minimap prepoznaje kao očitavanja mapirana na "početak" i "kraj" reference.

Početna pozicija na očitavanju	Krajnja pozicija na očitavanju	Početna pozicija na referenci	Krajnja pozicija na referenci	Duljina reference
46	11545	4630203	4641628	4641652
11575	19354	6	7676	4641652

**Tablica 4.4:** Mapiranje očitavanja na "rubove" cirkularnog kromosoma

Ta se očitavanja čine kao kimerna očitavanja mapirana na krajnje udaljene dijelove reference, ali zapravo je to jedinstveno očitavanje koje se nalazi na istom dijelu cirkularnog kromosoma i zato se ne smatra kimernim očitanjem nego regularnim.

Treći slučaj su očitavanja koja djeluju kao kimerna iako zapravo nisu. To su očitavanja s više mapiranja koja gledajući samo pozicije mapiranja na očitavanju djeluju kao kimerna. S druge strane, kada se pogledaju pozicije mapiranja na referenci, uz uvjet da su mapiranja na istom relativnom slijedu, vidljivo je da je praznina između mapiranja na očitavanju otprilike iste duljine kao praznina između mapiranja na referenci. Iz toga se može zaključiti da takvo očitavanje nije zapravo kimerno nego regularno koje iz nekog razloga, najčešće grešaka u sekvenciranju, ima jedan dio koji se lošije mapira na referencu.

I na kraju, postoje očitavanja koja kao rezultat minimap2-ovog algoritma imaju više mapiranja koja djeluju kao kimerna, ali među tim mapiranjima postoji najveće mapiranje koje sadrži barem 90% duljine cijelog očitavanja. Iako se takvo očitavanje u početku čini kao kimerno, ono je zapravo regularno očitavanje s manjom greškom zbog koje se sitan dio očitavanja mapira na drugi dio reference.

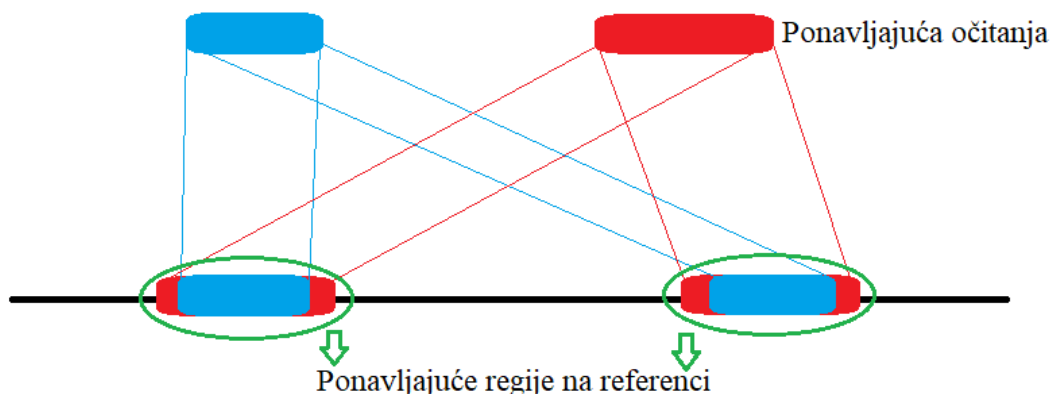
## 4.4. Detekcija ponavljajućih regija na referenci

Budući da su ponavljajuće regije, kao što im i ime kaže, ponavljajući dijelovi na referenci, za očekivati je da je na jednu regiju potencijalno moguće mapirati više različitih očitavanja.

Iz algoritma minimap2 vidljivo je da on, ako postoje, pronalazi višestruka mapiranja pojedinih očitavanja. Svako očitavanje koje u sebi sadrži regiju koja zapravo pripada ponavljajućoj regiji na referenci će biti mapirana na svaku od tih regija koja je vrlo slična ili ista onom dijelu reference iz kojeg je očitavanje zapravo uzeto.

Takva očitavanja smo već detektirali i klasificirali kao ponavljajuća očitavanja i upravo među njihovim mapiranjima tražimo ponavljajuće regije reference.

U procesu detektiranja ponavljajućih regija na referenci, među ponavljajućim očitavanjima tražimo višestruka (barem dvostruka) mapiranja različitih očitavanja na iste regije reference te one regije za koje postoje takva mapiranja proglašavamo ponavljajućim.



**Slika 4.6:** Primjer ponavljajuće regije na referenci detektirane pomoću dva ponavljajuća očitavanja koja se mapiraju na isto područje

Poznavajući ponavljajuće regije na referenci i očitavanja koja se tamo nalaze te ki-merna očitavanja, možemo prilikom izgradnje genoma obratiti posebnu pažnju na ta očitavanja i te regije kako bismo što uspješnije izbjegli zamke zbog kojih bi sastavljeni genom mogao izgledati vidno drugačije od stvarnog genoma.

Znajući sve to, ulazne FASTA podatke prilagođavamo onome što je detektirano.

## 4.5. Izlaz programa i sastavljanje genoma

Program datoteku s očitanjima u FASTA ili FASTQ formatu, koja je ulazni podatak programa, parsira uz pomoć biblioteke "bioparser". Iterirajući po vektoru očitavanja dobivenom parsiranjem datoteke pomoću jedinstveno identificirajućeg imena očitavanja pretražuju se kreirani setovi kimernih i ponavljajućih očitavanja te se ovisno o tome je li očitavanje detektirano i označeno kao kimerno, ponavljajuće, ponavljajuće sadržano ili je regularno očitavanje ono obrađuje u skladu s klasifikacijom. Očitavanja se uz informacije o prekidima za kimerna te informacije o ponavljajućim regijama odvajaju u tri FASTA datoteke. Opis obrade i zapisivanja očitavanja ovisno o vrsti očitavanja je u nastavku.

U prvoj FASTA datoteci nalaze se sva kimerna očitavanja. Svako očitavanje zapisano je sukladno standardnom FASTA formatu uz neke dodatke. U prvoj liniji se osim imena očitavanja nalaze i pozicije na očitavanju koje su detektirane kao prekidi, odnosno pozicije na kojima su očitavanja lažno spojena u jedno očitavanje. Datoteka se sprema u direktorij iz kojeg je program pokrenut pod imenom "chimeric\_reads.fasta".

U drugoj FASTA datoteci nalaze se ona ponavljajuća očitavanja koja se u cijelosti nalaze u nekoj ponavljajućoj regiji na referenci. Ta očitavanja ne mogu premostiti ponavljajuću regiju i zbog toga stvaraju probleme prilikom sastavljanja genoma. I ta očitavanja osim uobičajenih informacija sadržanih u FASTA formatu imaju pridodane neke dodatne informacije. U prvoj liniji se osim imena očitavanja nalaze i pozicije ponavljajućih regija toga očitavanja. Datoteka se sprema u direktorij iz kojeg je program pokrenut pod imenom "repeating\_reads.fasta".

Treću FASTA datoteku čine sva regularna očitavanja, kimerna očitavanja razdvojena po pozicijama koje su detektirane kao prekidi te ponavljajuća očitavanja koja nisu sadržana u nekoj ponavljajućoj regiji. Prilikom rezanja kimernih očitavanja pazilo se na ponavljajuće regije kimernih očitavanja te su u datoteku uvrštena samo ona izrezana očitavanja koja nisu ponavljajuća ili ako su ponavljajuća se ne nalaze u cijelosti u nekoj ponavljajućoj regiji na referenci. Ponavljajuća očitavanja uvrštena u datoteku uz ime sadrže i anotaciju s informacijama o ponavljajućim regijama unutar očitavanja. Datoteka s regularnim očitavanjima koristi se kao ulaz u alat za sastavljanje genoma Ra, a sprema se u direktorij iz kojeg je program pokrenut pod imenom "cleaned\_sequences.fasta".

Ra je assembler koji na temelju ulaznih očitavanja u FASTA/FASTQ formatu generira set kontiga visoke točnosti u FASTA formatu (Vaser i Šikić, 2019).

## 5. Implementacija i korištenje

Cjelokupna implementacija pisana je programskim jezikom C++. Kao ulaz program očekuje datoteku u PAF formatu koja sadrži rezultat mapiranja očitavanja na referencu te datoteku u FASTA ili FASTQ formatu koja sadrži početni skup očitavanja.

Osim programskog dijela pisanog u C++-u korištene su neke vanjske biblioteke te alati pomoću kojih je izrađeno i testirano programsko rješenje. Neki od tih alata i biblioteka su već spomenuti, a u nastavku je ukupan pregled svih korištenih tehnika te upute za korištenje svih potrebnih komponenti i samoga rješenja.

Za pretprocesiranje podataka korištene su dvije komponente. Prva komponenta je gore spomenuti alat `minimap2` koji uzimajući u obzir preddefinirane argumente mapira očitavanja na referencu te generira podatke o mapiranjima u PAF formatu.

Ulaz u programski dio rezultat je mapiranja u PAF formatu. Ti su podaci potom parsirani još jednom vanjskom komponentnom, `bioparserom`. `Bioparser` je višenamjenski parser bioinformatičkih formata. Osim za parsiranje PAF formata, korišten je i za parsiranje FASTA datoteke kako bi se među svim očitanjima iz FASTA datoteke izdvojila kimerna i ponavljajuća očitavanja od onih regularnih. Dobivene PAF i FASTA datoteke se parsirane spremaju u vektore `unique_ptr` koji pokazuju na strukture u kojima su pohranjene sve potrebne informacije poput imena očitavanja i reference, duljina očitavanja i reference, pozicija mapiranja, relativnih sljedova i slično te se podaci spremaju u tom vektoru koriste za detekcije i na kraju generiranja kimernih, ponavljajućih i regularnih FASTA datoteka.

Nakon detekcije i odvajanja kimernih, ponavljajućih i regularnih očitavanja, FASTA datoteka s regularnim očitanjima prosljeđena je alatu `Ra` koji iz pročišćenih očitavanja `overlap-layout-konsensus` pristupom `de novo` gradi genom te generira set kontiga u FASTA formatu (Vaser i Šikić, 2019).

Naposlijetku, za potrebe testiranja korišten je alat `QUAST-LG`. `QUAST-LG` je alat koji evaluira sastavljenost (engl. *assembly*) genoma koristeći razne metrike. Uspoređujući set kontiga izgeneriran alatom `Ra` s referentom bazom generira niz statistika i grafova koji pokazuju uspješnost sastavljenosti genoma u usporedbi s referencom.

Ovo su upute za korištenje svih navedenih komponenata u operacijskom sustavu Linux.

Za pripremu alata minimap2 za korištenje potrebni su ovi koraci:

```
$ git clone https://github.com/lh3/minimap2
$ cd minimap2 && make
```

Nakon toga je minimap2 spreman za korištenje. Minimap2 kao ulaz očekuje datoteku s očitajima i datoteku s referencom. Poziva se na ovaj način:

```
$ minimap2 reference.fasta reads.fasta > mapping.paf
```

Nakon toga slijedi detekcija kimernih i ponavljajućih očitajna. Kompletna implementacija nalazi se na platformi Github te se priprema za rad na ovaj način:

```
$ git clone --recursive
https://github.com/lbcb-edu/BSc-thesis-18-19.git
detection
$ cd detection
$ mkdir build
$ cd build
$ cmake -DCMAKE_BUILD_TYPE=Release ..
$ make
```

Detaljne upute za instalaciju, korištenje i te podaci su dostupni na <https://github.com/lbcb-edu/BSc-thesis-18-19/tree/sbakic>, a ovo su upute za pozivanje programa:

```
$ detection mapping.paf sequences.fasta
```

Za instalaciju alata Ra potrebno je izvršiti ove korake:

```
$ git clone --recursive
https://github.com/rvaser/ra.git ra
$ cd ra
$ mkdir build
$ cd build
$ cmake -DCMAKE_BUILD_TYPE=Release ..
$ make
```

Potom se Ra poziva naredbom:

```
$ ra -x {ont, pb} reads.fasta > contigs.fasta
```

## 6. Testiranje i rezultati

Rezultati su testirani gore spomenutim alatom QUAST-LG. Ono što je zapravo testirano jest kako se kontizi koje Ra izgradi nakon uklanjanja kimernih i ponavljajućih očitavanja podudaraju s postojećom referencom. Gledaju se brojne metrike pomoću kojih se ocjenjuje uspješnost sastavljenosti genoma te se generira opći izvještaj o uspješnosti sastavljenosti genoma (Gurevich et al., 2018a).

Kako bi bilo što lakše pratiti rezultate testiranja u nastavku je kratki pregled najvažnijih dijelova QUAST-ovog izvještaja. Izvještaj možemo podijeliti u tri tematska dijela.

Na samom početku nalaze se osnovne informacije o kontizima, broj kontiga, duljinama, udjelima G i C nukleotida u ukupnoj duljina kontiga u odnosu na udjele G i C nukleotida u ukupnoj duljini reference te duljinama i brojevima kontiga koji zadovoljavaju 50, odnosno 75 posto ukupne duljine.

U ovom dijelu izvješća nalaze se i informacije o sastavljenosti genoma, uspješnosti poravnanja kroz informacije o udjelima podudaranja, nepodudaranja i praznina na 100kbp i slični podaci vezani uz poravnate kontige. Iz tih informacije je moguće odrediti mjeru *identity* koja pokazuje udio genoma koji ima ispravno podudaranje nukleotida po poziciji.

Nakon tih podataka slijede informacije o pogrešno sastavljenim dijelovima (engl. *misassemblies*). Taj dio izvješća posebno je zanimljiv jer su pogrešno sastavljeni dijelovi greške najčešće uzrokovane ponavljajućim regijama i kimernim očitanjima, odnosno greškama u sekvenciranju. Pogrešno sastavljanja izuzetno su problematična u vidu korištenja sastavljene reference za daljnja istraživanja (Gurevich et al., 2018b).



QUAST-LG neku poziciju proglašava pozicijom pogrešnog sastavljanja ako zadovoljava jedan od ovih kriterija (Gurevich et al., 2018a):

- lijevo bočno očitavanje se poravnava više od 1 kbp (kbp=kilobase pairs) od desnog bočnog očitavanja na referenci (relokacija)
- bočna očitavanja se preklapaju više od 1kbp (relokacija)
- bočna očitavanja se poravnavaju na različite kromosome (translokacija) ili slijedove (inverzija)

QUAST-LG u izvješću pruža podatke o broju pogrešnih sastavljanja detaljizirano podijeljenih na gore definirane relokacije, translokacije i inverzije te podatke o kontizima koji ih sadržavaju.

U zadnjem dijelu su informacije o kontizima koji nisu uspješno poravnati, bilo djelomično ili u potpunosti.

Nakon tabličnog prikaza uspješnosti sastavljenosti genoma slijede grafički prikazi istih vrijednosti od kojih se kao najzanimljiviji izdvajaju slijedeći grafovi:

- graf udjela G i C nukleotida u referenci i u kontizima
- broj pogrešnih sastavljanja po relokacijama, translokacijama i inverzijama
- Feature-Response krivulja koja prikazuje ukupan broj poravnatih baza u kontigu podijeljen s ukupnom duljinom reference ovisno o broju značajki kontiga
- graf kumulativne duljine reference i kumulativne duljine poravnatih kontiga

Grafički prikazi QUAST-LG-ovom izvještaju se baziraju na gore opisanim mjerama. Ipak, treba pojasniti mjeru koja se ne pojavljuje eksplicitno u tabličnom dijelu izvješća, ali pogled na nju može dati uvid u kvalitetu sastavljanja. FRCurve(engl. Feature Response Curve) je mjera koja pokazuje kako rezultat sastavljanja odgovara na skup značajki. U konkretnom slučaju QUAST-LG-ovog izvješća pokazuje udio točno poravnatih baza u kontizima obzirom na broj značajki u kontigu (Gurevich et al., 2018a).

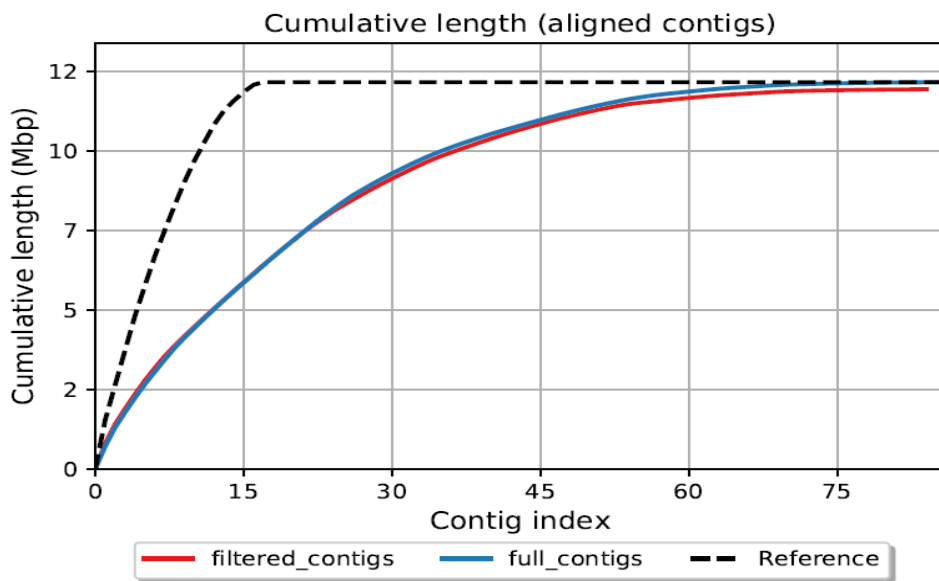
Testovi su rađeni na tri različita skupa podataka, *Saccharomyces Cerevisiae*, *Escherichia Coli* i *Klebsiella Pneumoniae*.

Prvi test rađen je na skupu podataka *Saccharomyces Cerevisiae*. Nakon detekcije kimernih i ponavljajućih očitavanja, 1164 kimerna očitavanja su izrezana po točkama prekida, 2569 ponavljajućih očitavanja koja nisu sadržana niti u jednoj ponavljajućoj regiji na referenci je anotirano po pozicijama ponavljajućih regija unutar očitavanja te je 3555 ponavljajućih očitavanja koja se u cijelosti sadržana u nekoj ponavljajućoj regiji izbačeno iz ulaznog skupa očitavanja. Sastavljena su dva genoma, jedan iz očitavanja dobivenih nakon detekcije kimernih i ponavljajućih očitavanja te jedan iz svih očitavanja bez detekcije te su neki od važnijih rezultata i grafova prikazani.

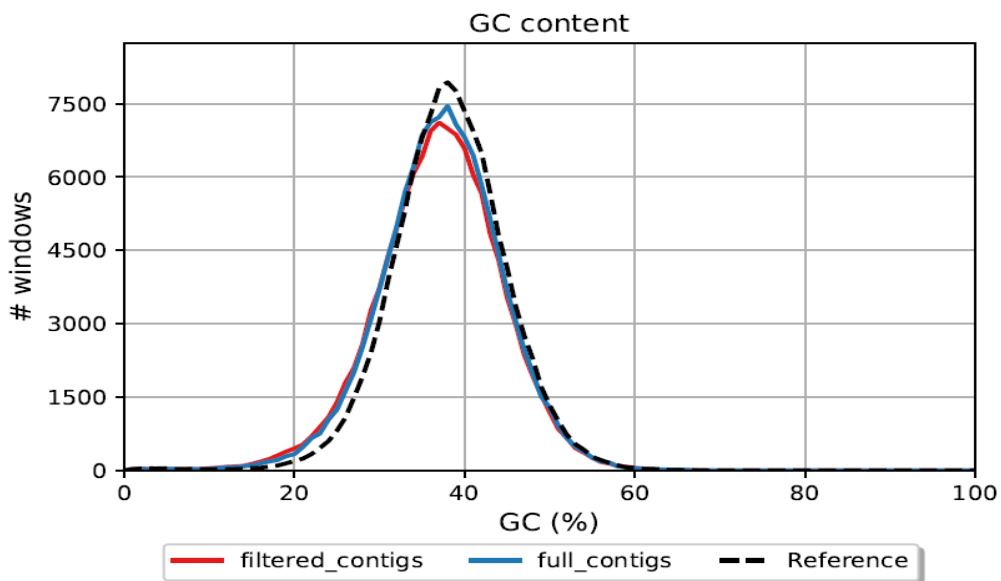
	Filtrirana očitavanja	Sva očitavanja
broj kontiga	68	69
broj misassembly-ja	16	17
sastavljenost genoma(%)	93.47	95.84
broj neporavnatih kontiga	0	0
identity(%)	98.8	98.9
duljina sastavljenog genoma	11 992 503	12 197 762
duljina ukupnog poravnanja	11 930 259	12 152 908
duljina reference		12 157 105

**Tablica 6.1:** Rezultati testiranja za *Saccharomyces Cerevisiae*

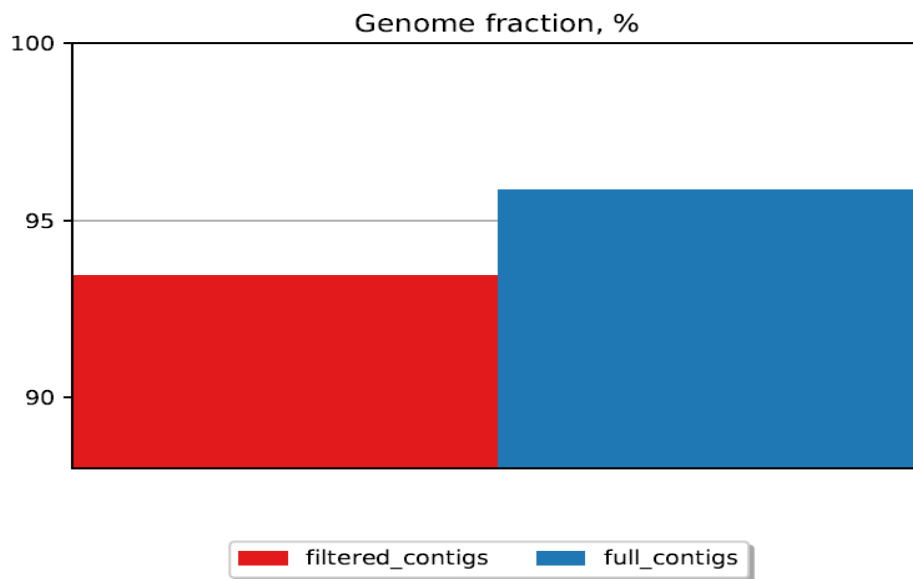
Vidljivo je da je uz manji broj kontiga smanjen i broj pogrešno sastavljenih dijelova, grešaka uzrokovanih ponavljajućim i kimernim očitanjima. Sastavljenost genoma malo je manja nego u slučaju sastavljanja svih očitavanja, ali to je potencijalno zbog razlike u pogrešnim sastavljanjima, a identity je približno isto uz malo kraću ukupnu duljinu genoma sastavljenog iz filtriranih očitavanja.



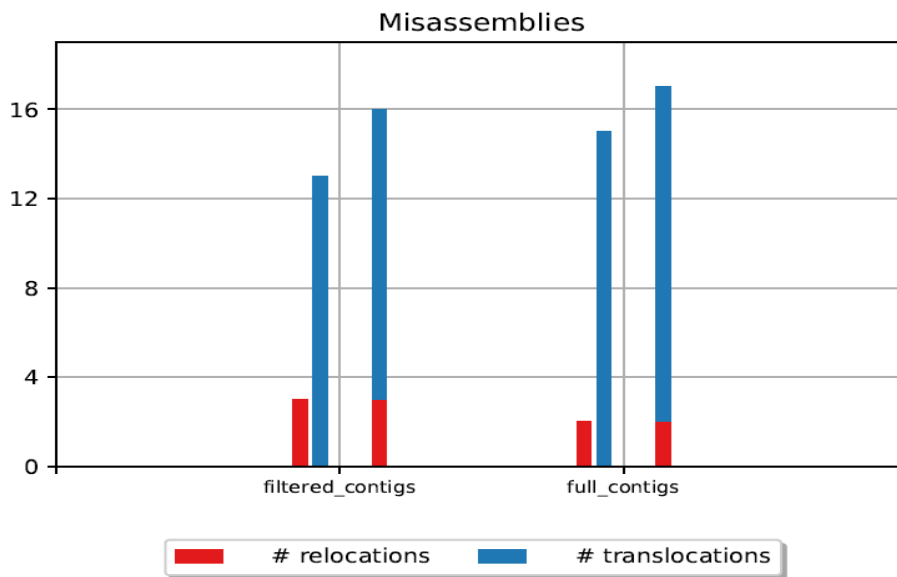
**Slika 6.1:** Odnos kumulativnih duljina kontiga dobivenih iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Saccharomyces Cerevisiae*



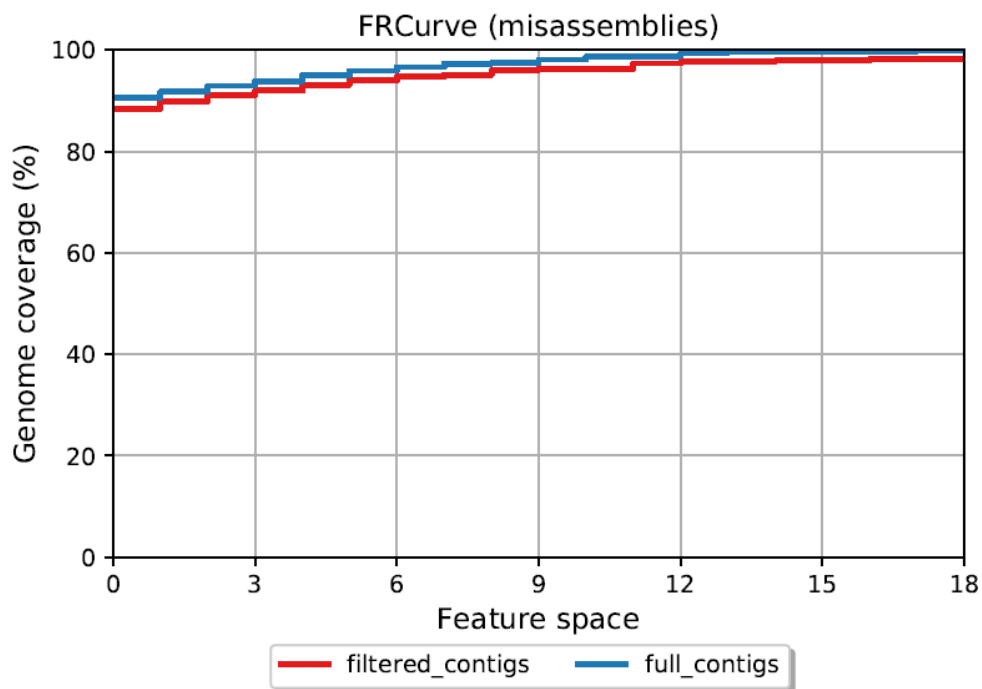
**Slika 6.2:** Odnos udjela G i C nukleotida u u kontizima dobivenim iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Saccharomyces Cerevisiae*



**Slika 6.3:** Odnos sastavljenosti genoma kontiga dobivenih iz filtriranih i nefiltriranih očitavanja za *Saccharomyces Cerevisiae*



**Slika 6.4:** Odnos broja pogrešnih sastavljanja za kontige dobivene iz filtriranih i nefiltriranih očitavanja za *Saccharomyces Cerevisiae*



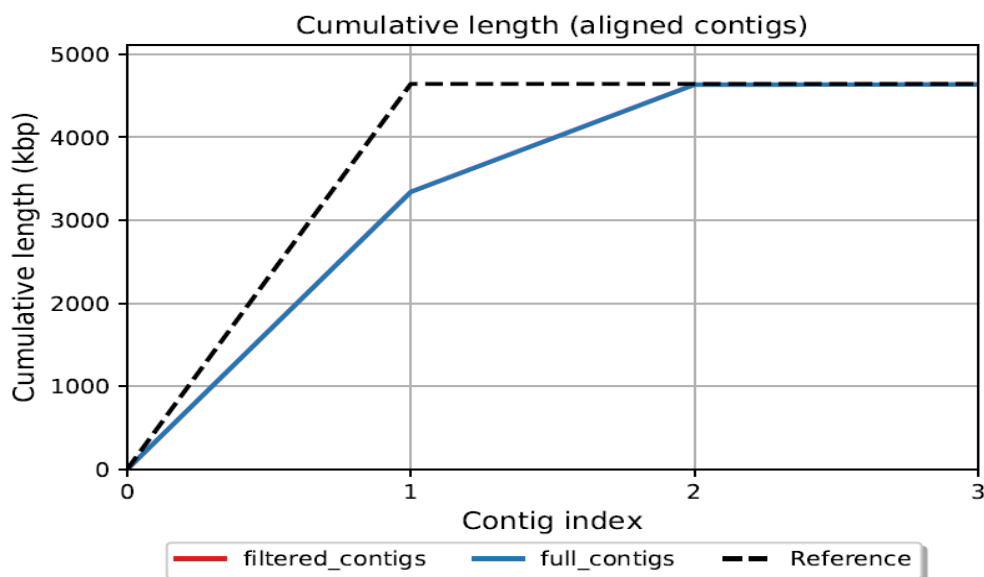
**Slika 6.5:** Odnos Feature Response-a za kontige dobivene iz filtriranih i nefiltriranih očitavanja za *Saccharomyces Cerevisiae*

I iz grafičkih podataka je vidljivo da je uz manji broj kontiga smanjen broj pogrešno sastavljenih dijelova. Ipak, kao problem ostaje malo slabija sastavljenost genoma u smislu duljine te sporija konvergencija potpunoj pokrivenosti genoma vidljiva na grafu Feature Response krivulje.

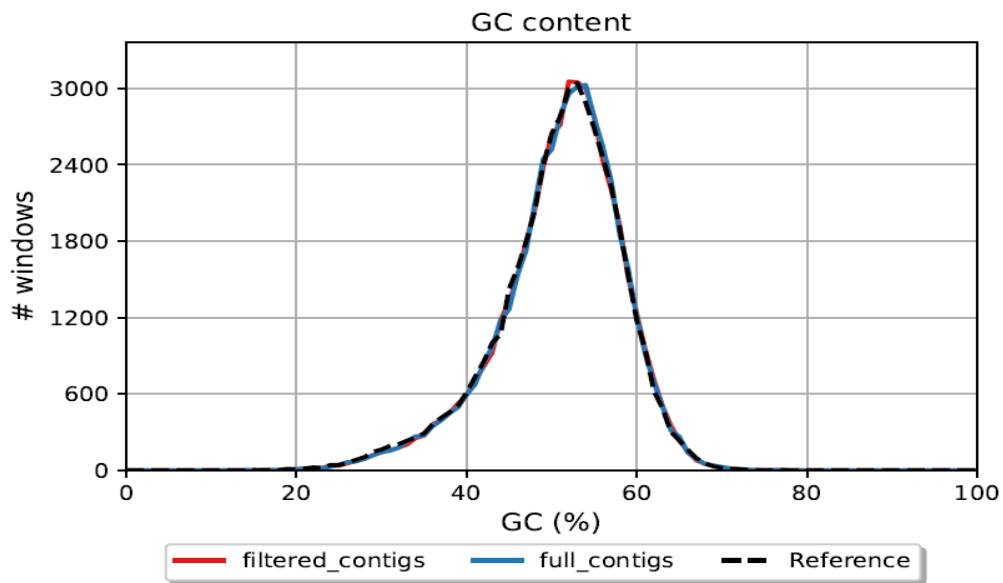
Na isti način je proveden test na skupu podataka *Escherichia Coli*. Iz skupa svih očitavanja izdvojeno je 71 očitavanje kao kimerno te je na točkama prekida izrezano na nova regularna očitavanja, 9 očitavanja prepoznato je kao ponavljajuće očitavanje sadržano u nekoj ponavljajućoj regiji na referenci i kao takvo je izbačeno iz konačnog skupa očitavanja koji je ulaz alatu Ra, a 57 očitavanja je klasificirano kao ponavljajuće očitavanje koje nije u cijelosti sadržano u nekoj ponavljajućoj regiji na referenci te je anotirano i prosljeđeno sastavljanju genoma. Nakon testiranja alatom QUAST-LG dobiveni su ovi rezultati:

	Filtrirana očitavanja	Sva očitavanja
broj kontiga	1	1
broj misassembly-ja	2	2
sastavljenost genoma(%)	99.988	99.957
broj nepravilnih kontiga	0	0
identity(%)	99.5	99.5
duljina sastavljenog genoma	4 635 596	4 630 485
duljina ukupnog poravnanja	4 635 595	4 630 463
duljina reference	4 641 652	

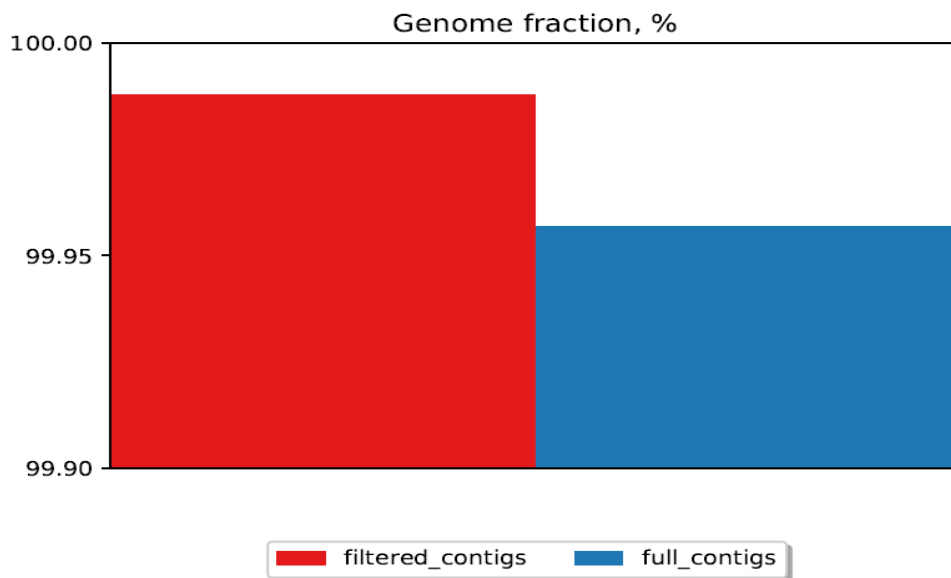
**Tablica 6.2:** Rezultati testiranja za *Escherichia Coli*



**Slika 6.6:** Odnos kumulativnih duljina kontiga dobivenih iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Escherichia Coli*



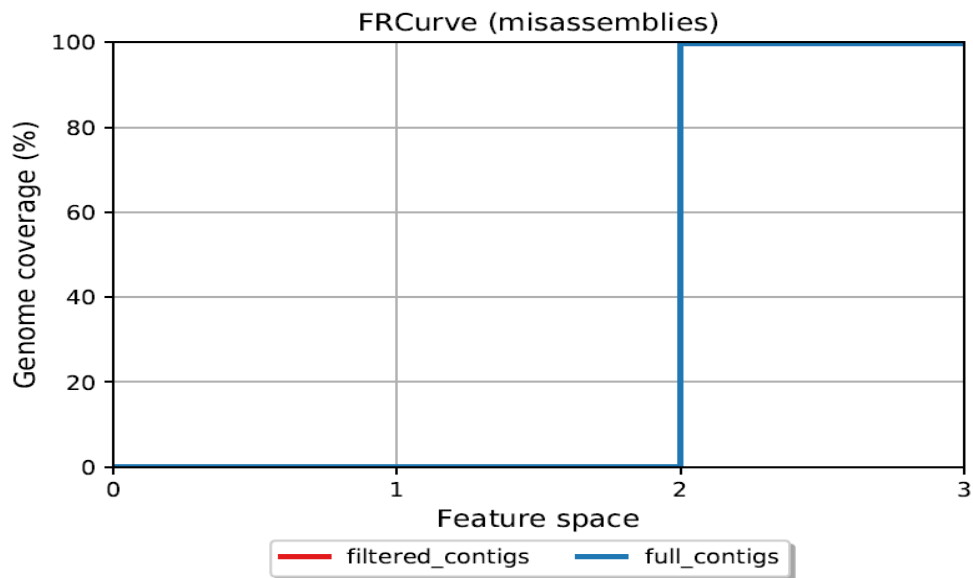
**Slika 6.7:** Odnos udjela G i C nukleotida u u kontizima dobivenim iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Escherichia Coli*



**Slika 6.8:** Odnos sastavljenosti genoma kontiga dobivenih iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Escherichia Coli*



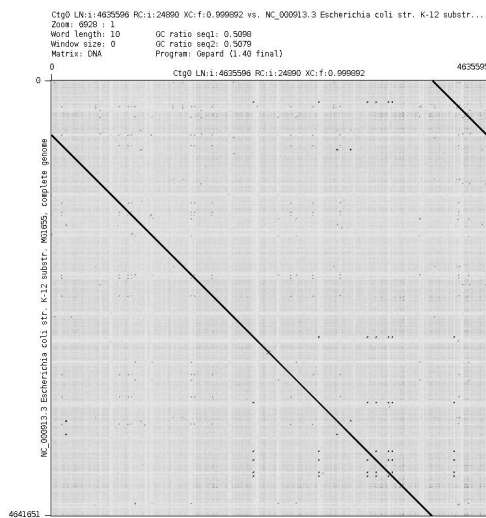
**Slika 6.9:** Odnos broja pogrešnih sastavljanja za kontige dobivene iz filtriranih i nefiltriranih očitavanja za *Escherichia Coli*



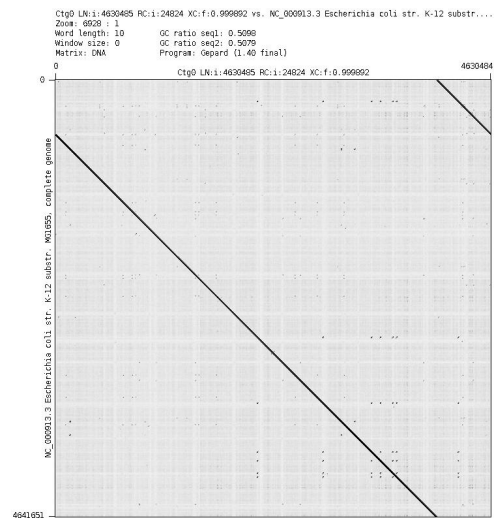
**Slika 6.10:** Odnos Feature Response-a za kontige dobivene iz filtriranih i nefiltriranih očitavanja za *Escherichia Coli*



Rezultati ovoga testiranja pokazuju da se broj kontiga prije i nakon filtriranja nije promijenio. Ipak, treba imati na umu da je *Escherichia Coli* relativno "malena" (samo 4641652 baza) te da je prepoznat malen broj kimernih (71) i ponavljajućih (66) očitavanja. Sastavljenost genoma malo je bolja nakon filtracije, a tvrdnja da je genom zaista dobro sastavljen, usprkos dva pronađena pogrešna sastavljanja može se potvrditi uspoređivanjem sastavljenog genoma i postojeće reference u alatu gepard. Gepard kao ulaz očekuje dvije datoteke u FASTA ili FASTQ formatu sa sastavljenim genomima pomoću kojih iscrtava točkasti graf s pozicijama jednog genoma na x-osi te pozicijama drugog genoma na y-osi. Dobro podudaranje zapravo znači poklapanje na istoj poziciji što je praktično pravac nagnut pod kutom od 45° u odnosu na x-os (Krumisiek et al., 2007).



(a) Poravnanje genoma izgrađenog iz filtriranih očitavanja na referencu



(b) Poravnanje genoma izgrađenog iz svih očitavanja na referencu

Uz jednake iznose identity-ja od 99.5% iz slika je vidljivo da se sastavljeni genomi i postojeća referenca gotovo savršeno poklapaju stoga se da zaključiti da su prepoznata pogrešna sastavljanja neke manje greške koje nisu bitno utjecale na krajnji izgled genoma.

Posljednji test rađen je na skupu podataka *Klebsiella Pneumoniae*. Iz početnog skupa podataka je detektirano i po pozicijama prekida izrezano 5345 kimernih očitavanja. 721 ponavljajuće očitavanje se nije nalazilo niti u jednoj ponavljajućoj regiji na referenci te je anotirano i dodano u skup očitavanja za sastavljanje, a 272 očitavanja su bila ponavljajuća i u potpunosti sadržana u nekoj ponavljajućoj regiji na referenci. Testiranje je dalo ovakve rezultate:

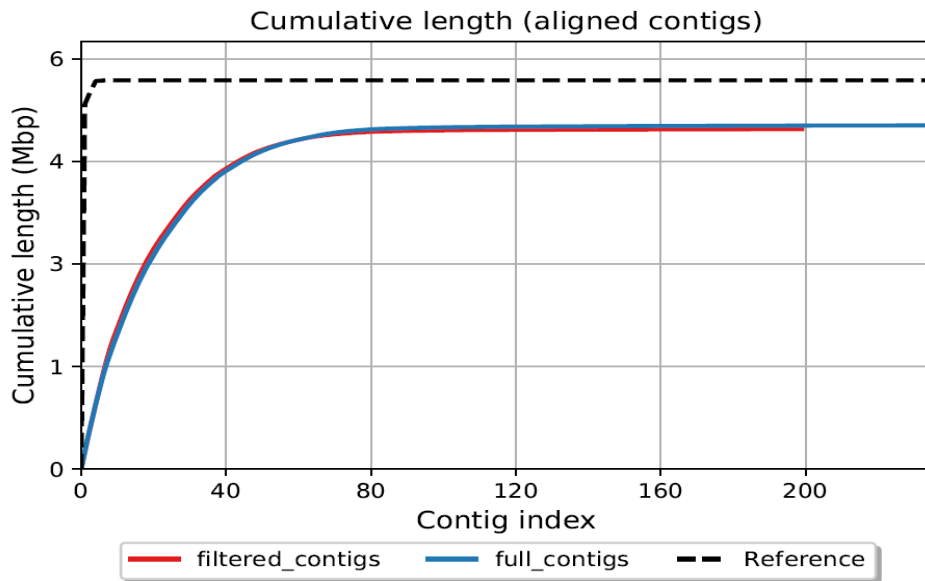
	Filtrirana očitavanja	Sva očitavanja
broj kontiga	16	9
broj misassembly-ja	71	85
sastavljenost genoma(%)	87.29	86.94
broj neporavnatih kontiga	0	0
identity(%)	98.687	98.686
duljina sastavljenog genoma	5 481 106	5 559 831
duljina ukupnog poravnanja	4 968 356	5 024 577
duljina reference	5 682 322	

**Tablica 6.3:** Rezultati testiranja za *Klebsiella Pneumoniae*

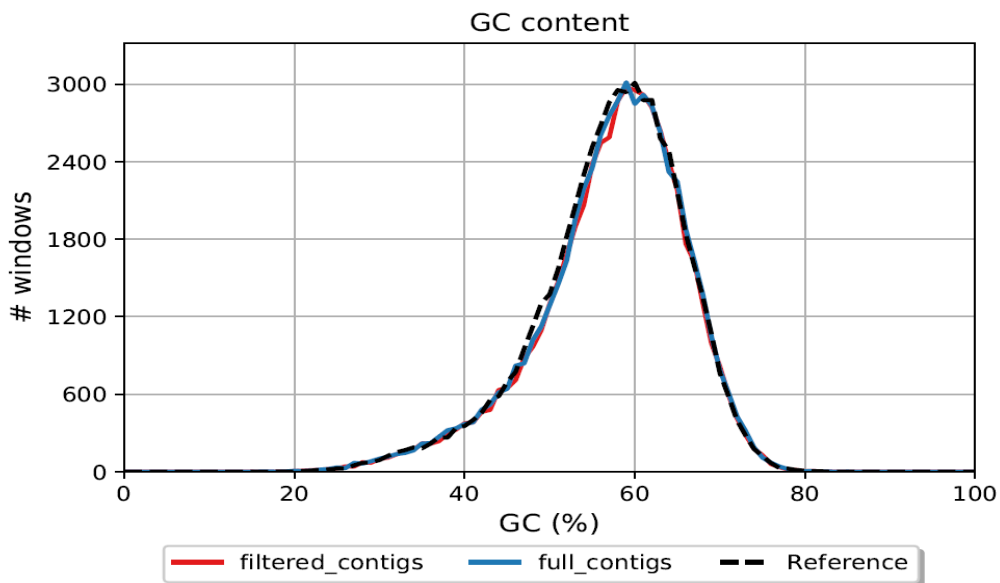
Skup podataka za *Klebsiella Pneumoniae* dao je rezultate koji nakon filtriranja podataka imaju više kontiga te malu razliku u sastavljenosti genoma. S druge strane, broj pogrešnih sastavljanja nakon filtriranja vidljivo je manji u odnosu na originalni skup očitavanja.

Postoji više potencijalnih razloga za ovakve rezultate. Jedan od mogućih razloga slabija je kvaliteta referentne baze podataka ili baze podataka očitavanja zbog čega prilikom mapiranja očitavanja na referencu dolazi do, u prosjeku, lošijih mapiranja i zbog toga se više očitavanja detektira kao kimerna očitavanja što kasnije može stvarati probleme prilikom sastavljanja genoma.

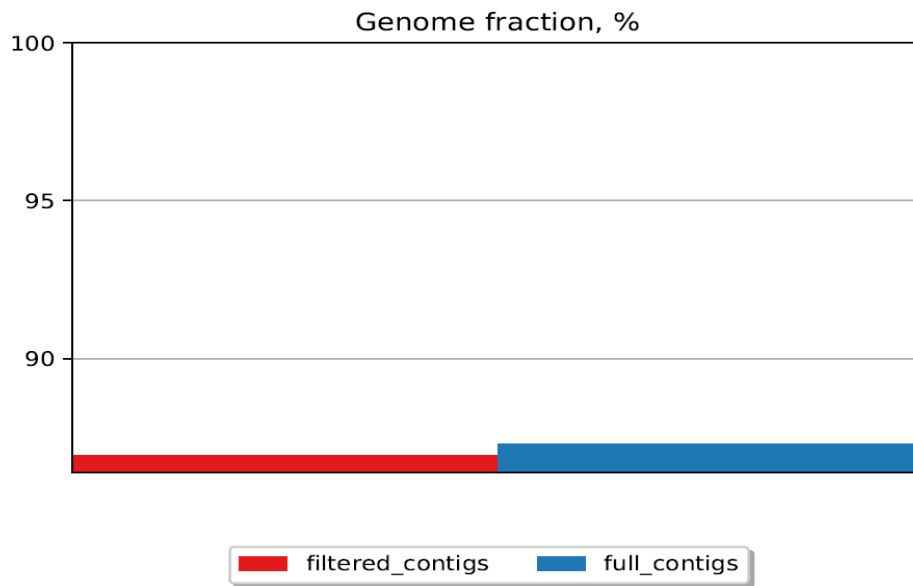
Ipak, ono što je pozitivan rezultat ovoga testa jest činjenica da je, usprkos povećanju broja kontiga, smanjen broj pogrešno sastavljenih dijelova u izgrađenom genomu.



**Slika 6.12:** Odnos kumulativnih duljina kontiga dobivenih iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Klebsiella Pneumoniae*



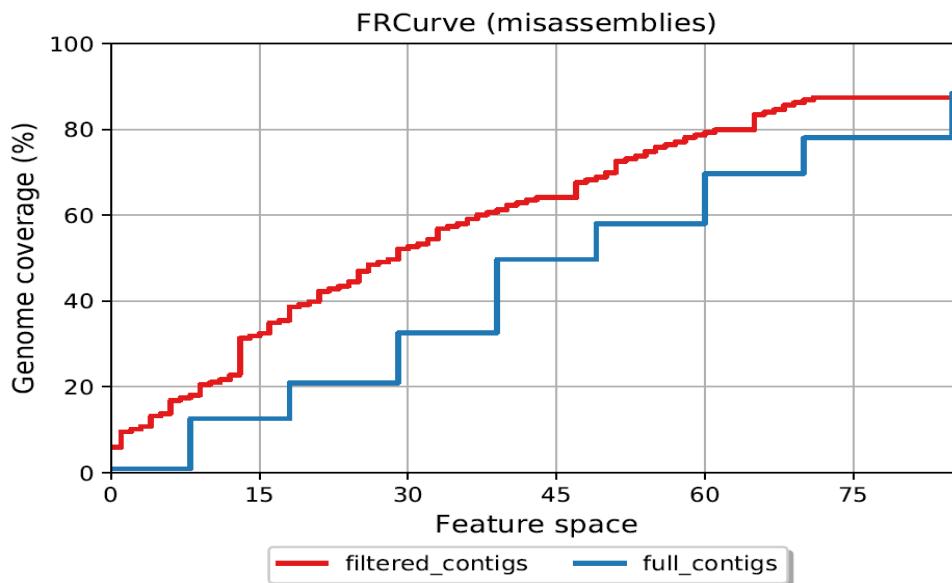
**Slika 6.13:** Odnos udjela G i C nukleotida u u kontizima dobivenim iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Klebsiella Pneumoniae*



**Slika 6.14:** Odnos sastavljenosti genoma kontiga dobivenih iz filtriranih i nefiltriranih očitavanja te usporedba s referencom za *Klebsiella Pneumoniae*



**Slika 6.15:** Odnos broja pogrešnih sastavljanja za kontige dobivene iz filtriranih i nefiltriranih očitavanja za *Klebsiella Pneumoniae*



**Slika 6.16:** Odnos Feature Response-a za kontige dobivene iz filtriranih i nefiltriranih očitavanja za *Klebsiella Pneumoniae*

Obzirom da je većina grafičkih prikaza već objašnjena, kao posebno zanimljiv ostaje graf Feature Response krivulje. Naime, na grafu je vidljivo da za filtrirana očitavanja prekrivenost genoma ima puno bolji odziv u odnosu na kontige sastavljene iz svih očitavanja za isti broj značajki.

## 6.1. Daljnji rad

Kroz ovih nekoliko testova pokazano je da za de novo sastavljanje genoma postoje određeni benefiti u pronalasku i anotaciji kimernih i ponavljajućih očitavanja.

Pokazalo se, doduše, da postoje i određeni problemi na koje ovakav način detekcije nailazi ako su podaci za mapiranje lošije kvalitete.

Kako bi detekcija što više pomogla prilikom sastavljanja genoma, bilo bi idealno osim podataka o mapiranjima očitavanja na postojeću referencu koristiti još neke informacije o samim očitanjima i na taj način još kvalitetnije selektirati stvarna kimerna i ponavljajuća očitavanja među ostalim očitanjima.

## 7. Zaključak

Bioinformatika kao interdisciplinarno područje okuplja širok spektar znanja i stručnjaka koji provode kompleksna istraživanja s ciljem pronalaska što točnijih i konkretnijih informacija vezanih uz bioznanost. Bioznanost kao vrlo osjetljivo područje zahtijeva podatke koji su vrlo visoke kvalitete i preciznosti.

S padom cijena istraživanja počeo je ubrzani razvoj s novim metodama sekvenciranja. To je omogućilo veći broj istraživanja, ali je uzrokovalo i veću stopu pogrešaka u procesu sekvenciranja. Problemi odabira odgovarajućih očitavanja na ponavljajućim područjima genoma te izbjegavanje pogrešnog sastavljanja genoma uslijed kreiranja "umjetnih" očitavanja, problemi su koji nisu izbjegnuti u novoj generaciji sekvenciranja. Takvi problemi nameću se kao ugrožavajući faktor za kvalitetu sastavljenih genoma, a samim time i kvalitetu istraživanja kojima se podvrgavaju tako dobiveni genomi.

Kada bi se i takvi problemi uspješno rješavali, primjena rezultata istraživanja bila bi još veća u svakodnevici te je zato vrlo bitno razvijati bioinformatičke alate i algoritme u smjeru što veće preciznosti.

Upravo zato važno je raditi na alatima i metodama koje će efikasno i precizno detektirati kimerna i ponavljajuća očitavanja te koristiti znanja o njima kako bi novo sastavljeni genom bio što točniji i upotrebljiviji u daljnjim testiranjima i istraživanjima. Takvi alati u suradnji s alatima za mapiranja i sastavljanje genoma mogu dovesti do genoma visoke kvalitete i upotrebljivosti što se nameće kao vrlo poželjan ishod.

Ovaj rad pokazuje kako detekcija kimernih i ponavljajućih očitavanja daje rezultate prilikom sastavljanja genoma. Uz detekciju i kvalitetno iskorištavanje informacija o kimernim i ponavljajućim očitavanjima moguće je smanjiti broj pogrešno sastavljenih dijelova genoma te tako popraviti kvalitetu sastavljanja. Naravno, kako bi rezultati detekcije bili što precizniji potrebno je nadalje razvijati metode i testirati pristupe detekciji takvih očitavanja.

# LITERATURA

- atdbio. Next generation sequencing, 2011. URL <https://www.atdbio.com/content/58/Next-generation-sequencing>.
- Alexey Gurevich, Vladislav Saveliev, Nikolay Vyahhi, i Glenn Tesler. Quast 5.0.2 manual, 2018a. URL <http://quast.bioinf.spbau.ru/manual.html#sec3>.
- Alexey Gurevich, Vladislav Saveliev, Nikolay Vyahhi, i Glenn Tesler. Quast: quality assessment tool for genome assemblies, 2018b. URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3624806/>.
- Amy M. Hauth i Deborah A. Joseph. Beyond tandem repeats: complex pattern structures and distant regions of similarity, 2002. URL [https://academic.oup.com/bioinformatics/article/18/suppl\\_1/S31/231738](https://academic.oup.com/bioinformatics/article/18/suppl_1/S31/231738).
- Mile Šikić i Mirjana Domazet-Lošo. *Bioinformatika*. 2013. URL [https://www.fer.unizg.hr/\\_download/repository/bioinformatika\\_skripta\\_v1.2.pdf](https://www.fer.unizg.hr/_download/repository/bioinformatika_skripta_v1.2.pdf).
- J. Craig Venter Institute. Genome sequencing, 2003. URL <http://www.genomenetwork.org>.
- Jan Krumsiek, Roland Arnold, i Thomas Rattei. Gepard: a rapid and sensitive tool for creating dotplots on genome scale, 2007. URL <https://academic.oup.com/bioinformatics/article/23/8/1026/198110>.
- Arthur M. Lesk. *Bioinformatics*. 2013. URL <https://www.britannica.com/science/bioinformatics>.
- Heng Li. Minimap2: pairwise alignment for nucleotide sequences, 2018a. URL <https://github.com/lh3/minimap2>.

Heng Li. Paf: a pairwise mapping format, 2018b. URL <https://github.com/lh3/miniasm/blob/master/PAF.md>.

Robert Vaser i Mile Šikić. Yet another de novo genome assembler, 2019. URL <https://www.biorxiv.org/content/10.1101/656306v1.article-metrics>.



## **De novo sastavljanje genoma vođeno referencom**

### **Sažetak**

De novo sastavljanje genoma jedan je od najkompleksijih problema u bioinformatici. Brojni su problemi prilikom sastavljanja genoma, uzrokovani, kako biološkim specifičnostima, tako i greškama u procesu sekvenciranja. Probleme stvaraju takozvana ponavljajuća i kimerna očitavanja. U ovom radu prikazano je kako detektiranje kimernih i ponavljajućih očitavanja uz prikladna postupanja s njima pomaže ispravnom sastavljanju genoma. Prikazane su poteškoće na koje se nailazi prilikom sastavljanja genoma te kako uz postojeće alate minimap2 i Ra i ispravnu detekciju kimernih i ponavljajućih očitavanja što bolje sastaviti genom te rezultati testiranja sastavljenosti genoma. Osim toga, pojašnjeni su i formati korištenih podataka te temeljne značajke korištenih alata.

**Ključne riječi:** kimerno, onavljajuće, genom, minimap2, ra, FASTA, FASTQ, PAF, bioinformatika.

## **Reference-Guided de novo Genome Assembly**

### **Abstract**

De novo genome assembly is one of the most complex problems in Bioinformatics. There are many problems, either caused by biological specificities or mistakes made during sequencing process, that complicate de novo genome assembly process. The problems are mainly caused by chimeric and repeating reads. This thesis describes how detecting and annotating chimeric and repeating reads helps the assembly process. Typical issues that occur during the assembly process and how successful existing tools Ra and minimap2 with the detection of chimeric and repeating reads assembly the reads are shown. Besides that, thesis describes file formats used in this work and thorough features of used tools

**Keywords:** chimers, repeats, genome, minimap2, Ra, FASTA, FASTQ, PAF, Bioinformatics