

Asembliranje sekvenci

Seminarski rad u okviru kursa
Uvod u bioinformatiku
Matematički fakultet

Aleksandra Karadžić, Dragutin Ilić
karadzic.matf@gmail.com, dragutin_ilic@yahoo.com

18. maj 2017.

Sažetak

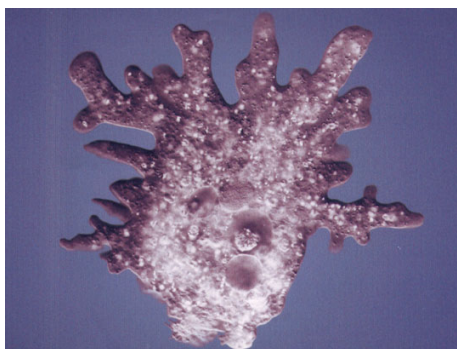
Sekvenciranje genoma predstavlja jedan od osnovnih zadataka u bioinformatiki. Kako još uvek nije dostupan način kojim bi smo mogli pročitati celu DNK sekvencu bioinforamtičari su posvetili dosta pažnje usavršavanju metoda za asembliranje sekvenci. U ovom radu je opisan jedan *de-novo* metod zasnovan na preklapanju ridova kratke dužine, način na koji se formiraju kontige kao i dobijanje konsenzus sekvence kontiga. Takođe prikazan je i rad programa PCAP.Solexa koji predstavlja implementaciju pomenutog metoda.

Sadržaj

1	Uvod	2
2	Opis algoritma	4
2.1	Konstrukcija niza super-reč	5
2.2	Izračunavanje preklapanja koristeći niz super-reč	6
2.3	Ubrzanje paralelizacijom	6
2.4	Konstrukcija grafa preklapanja	6
2.5	Generisanje fajla konsenzus sekvenci kontiga i ocene kvaliteta	7
3	Korišćenje programa PCAP.Solexa	7
3.1	Priprema, pokretanje i podaci korišćeni pri radu programa .	7
3.2	Rezultati rada programa	8
3.3	Potencijalni problemi	10
4	Zaključak	10
	Literatura	10
A	Dodatak	11
A.1	Ocena kvaliteta sekvenci Phred	11

1 Uvod

Sekvenciranje genoma, odnosno određivanje rasporeda nukleotida (nt) u genomu, je jedan od fundamentalnih zadataka u bioinformatici. Nalazi svoju primenu u mnogim poljima nauke, uključujući funkcionalnu genetiku, onkologiju, evolucionarnu biologiju, forenzičke nauke i mnoge druge. Dužina genoma varira u zavisnosti od organizma (dužina ljudskog genoma je oko 3 milijarde nt, a jedan od najdužih genoma pripada amorfnom jednoćelijskom organizmu *Ameoba dubia* koja je oko 200 puta duži, slika 1). Na slici 2 su prikazane veličine genoma različitih organizama.



Slika 1: Ameoba dubia

Genome	Base pairs	No. of Genes
Phi-X 174	5,386	10
<i>Nanoarchaeum equitans</i>	490,885	552
<i>E. coli</i>	4,639,221	4,377
<i>Saccharomyces cerevisiae</i>	12,495,682	5,800
<i>Drosophila melanogaster</i>	122,653,977	13,379
<i>Homo sapiens</i>	3.2×10^9	30,000
<i>Protopterus aethiopicus</i>	1.3×10^9	?
<i>Psilotum nudum</i>	2.5×10^{11}	?20-25,000
<i>Amoeba dubia</i>	6.7×10^{11}	?

Slika 2: Broj baznih parova, kao i broj gena koje imaju ćelije različitih organizama. Obeleženo žutim je veličina ljudskog genoma.

Problem koji se javlja je zapravo činjenica da još uvek nisu razvijene tehnologije koje omogućavaju čitanje nukleotida u genomu od početka do kraja (slično kao čitanje stranica knjige sleva udesno). Trenutno zastupljeno rešenje je sekvenciranje manjih fragmenata DNK koji se nazivaju

ridovi (*eng. reads*). Uzima se mali uzorak tkiva ili krvi koji sadrži milione kopija DNK. Biohemijskim procesima se DNK razbija na fragmente, čijim sekvenciranjem se dobijaju ridovi. Ne zna se iz kog dela genoma je dobijen određeni rid, pa se koristi tehnika preklapanja ridova da bi se rekonstruisao genom. Ovaj ceo proces se naziva i **asembliranje genoma** ili **asembliranje sekvenci** (*eng. genome assembly, sequence assembly*).

Prvo sekvenciranje genoma je odrađeno 1977. godine od strane Frederika Sangera (*eng. Frederick Sanger*). U Sangerovom metodu dužina ridova je bila između 500 i 1000 nt. Većina programa zasnovana na ovoj tehnici (*eng. de novo assembly programs*) se bazira na “preklapanje-raspored-konsenzus strategiji” (*eng. overlap-layout-consensus strategy*), u kojoj se preklapanja između ridova izračunavaju brzim tehnikama upoređivanja, raspored kontiga (*eng. contigs*) se generiše pomoću preklapanja u opadajućem rasporedu kvaliteta, a konsenzus sekvenci kontiga se dobija brzim metodama višestrukog poravnavanja. Glavni razlog uspešnosti ovih programa jeste ta da je veličina ridova dovoljna za ustanovljavanje razlika između pravih i lažnih preklapanja. Ovakav način sekvenciranja pripada tehnologijama prve generacije kojim je proizvedena velika količina podataka DNK sekvenci. Iako sporije i skuplje tehnologije ove generacije se još uvek koriste tamo gde se zahteva preciznost.

Pre nego što uđemo u sledeću generaciju sekvenciranja (*eng. next-generation sequencing (NGS)*) navešćemo neke praktične probleme na koje nailazimo prilikom sekvenciranja genoma. Prvo, DNK se sastoji od dve niti, pa ne možemo znati iz koje od njih je rid izveden (nemamo informaciju da li da koristimo dobijen rid ili njegov obrnuti komplement prilikom asembliranja određene niti u genomu). Drugo, tehnologije koje se koriste nisu savršene, pa dobijeni ridovi često sadrže greške (čime se otežava preklapanje ridova). Treće, neki regioni genoma mogu da ostanu nepokriveni ridovima, čime je onemogućena rekonstrukcija celog genoma. Za otklanjanje ovih problema se koriste razne tehnike kao na primer: pristup razbijanja ridova (*eng. read breaking approach*) se koristi za problem nepokrivenosti genoma ridovima, asembliranje kontiga (neprkidnih delova genoma), a ne čitavih hromozoma, se vrši zbog navedenih problema, tehnika uklanjanja mehurova (*eng. bubble removal*) kojom se rešavaju ridovi sa greškom (lažni ridovi) itd.

Ulaskom u drugu generaciju (masovno paralelnih) sekvencijalnih tehnologija, počelo se sa proizvodnjom na milijarde kratkih ridova dužine od 50 do 150 nukleotida. Ogromni skupovi podataka vremenski i prostorno usložnjavaju izračunavanje preklapanja dinamičkim algoritmima. Ridovi male dužine otežavaju otkrivanje tačnih od lažnih preklapanja prilikom generisanja rasporeda kontiga. Kako bi se izbegle ove poteškoće, razvijeni su algoritmi za asembliranje zasnovani na detekciji preklapanja tačnih podudaranja fiksne dužine. Na osnovu ovakvih preklapanja formira se *de Bruijn*-ov graf, u kome svaki jedinstveni string dužine k (k -mer) koji se pojavljuje u ridu, predstavlja granu koja spaja čvorove označene sa $k-1$ -mer prefiksom i $k-1$ -mer sufixom. Pomoću ovog grafa i distribucije frekvencije stringova u ridovima razlikujemo stringove koji imaju grešku od onih koji nemaju. Stringovi sa greškama se ili ne koriste ili se prepravljaju i zatim koriste pri asembliranju. Iako se De Bruijnov graf može efikasno konstruisati i koristiti, problem je što se sekvence ridova duže od k ne koriste direktno u konstrukciji kontiga.

Kako bi koristili sekvenciranje ridova direktno, u ovom radu je predstavljen efikasan metod za izračunavanje preklapanja za kratke ridove.

Metod dozvoljava promašaje u preklapanjima, ali ne dodavanja i brisanja, koji se inače i javljaju ređe od promašaja. U grafu preklapanja, svaki čvor predstavlja rid, dok grana predstavlja preklapanje. U ovom grafu algoritam pronalazi jedinstvenu putanju ridova koja reprezentuje kontigu. Konsenzus sekvence svake kontige je dobijen izračunavanjem poravnanja višestrukih ridova koji nisu razdvojeni nukleotidima koji oni ne sadrže.

U nastavku će biti predstavljeni detalji algoritma kao i opis programa PCAP.Solexa u kome je algoritam implementiran.

2 Opis algoritma

Na početku ćemo opisati metod koji se koristi za izračunavanje preklapanja između ridova. Za realizaciju ovog metoda je korišćena struktura podataka koju ćemo zvati **“niz super-reč”** (*eng. superword array*), koja je ime dobila na sličan način kao i druga struktura podataka zvana **“sufiksni niz”** (*eng. suffix array*). Reč dužine w je string od w karaktera, a super-reč sa v reči je string sa $v \cdot w$ karaktera dobijena konkatencijom tih reču po redu u kome se javljaju. Reči u super-reči ćemo indeksirati pozicijama od 1 do v . Pozicije u rečima (indeksiranje kreće od jedinice) se dele u dve grupe: čekirane i nečekirane pozicije. Za dve reči dužine w kažemo da se podudaraju ukoliko imaju identične baze na svakoj čekiranoj poziciji tih reči. Na primer ako imamo dve reči dužine 15, **ACCATAC-CATAGCAC** i **ACTATTCCATAACAC**, i ako pretpostavimo da su nečekirane pozicije 3, 6 i 12 tada za ove dve reči možemo reći da se podudaraju. Dužina reči w se često namešta da bude 12 ili veća, kako bi broj čekiranih pozicija u reči bio 12, čime se osigurava da “tabela pogleda” (*eng. lookup table*) za sve stringove ove dužine (ako je reč duža od 12 brišemo sve nukleotide iz nečekiranih pozicija) može stati u glavnu memoriju. Broj reči v se bira tako da super-reč dužine $v \cdot w$ bude manja od dužine svakog rida (dužinu rida ćemo označavati sa r). Dužina super-reči se još naziva i minimalna dužina preklapanja. Adekvatan primer za ridove dužine 150 bi bio da postavimo da w na 15 i v na 6 čime dobijamo super-reči dužine 90. Rid dužine r ima $r - w \cdot v + 1$ super-reči na pozicijama 1, 2, ..., $r - w \cdot v + 1$. U našem primeru bi rid dužine 150 imao 61 super-reč na pozicijama 1, 2, ..., 61. Za dve super-reči kažemo da su identične ako im se sve reči podudaraju, odnosno svi nukleotidi na čekiranim pozicijama su im identični. Za jednu super-reč kažemo da je manja (veća) od druge ako string nukleotida na svakoj čekiranoj poziciji prve super-reči, u leksiografskom poredku, dolazi pre (posle) stringa nukleotida na čekiranim pozicijama druge super-reči. Svaki rid dobija jedinstveni nenegativni ceo broj koji se zove indeks rida. Pomoću njega i početne pozicije super-reči u ridu se izračunava jedinstveni indeks svake super-reči u svakom ridu. U algoritmu se koristi i funkcija koja pomoću datog indeksa super-reči efikasno računa poziciju te super-reči u indeksu i indeks rida. Sada imamo sve podatke kako bi definisali “niz super-reč” strukturu podataka. **“Niz super-reč”** za skup ridova je niz svih indeksa super-reči gde su super-reči sortirane od manje ka većoj.

Na slici 3 prikazano je osam super-reči u sortiranom poredku. Svaka super-reč se sastoji od četiri reči ($v=4$) dužine $w = 15$, gde je svaka čekirana pozicija označena bitom 1, a nečekirana 0. Poslednja pozicija svake reči je označena taraba znakom. Super-reči su grupisane tako da se u svakom bloku nalaze identične super-reči (prve tri, sledeće dve i poslednje tri). Vidimo da nukleotidi na nečekiranim pozicijama mogu biti

```

110110111110111110110111110111110110111110111110110111110111110111110111
# # # #
ATNAGCCCAGTTATCCTAGTCAGACTCAGGTTNCATCATTCCGANCAGACTGACCAG
ATGAGGCCAGTCATCCTTGTGAGACTTAGGTTACATCATTCCGAACAACTGANCAG
ATTAGNCCAGTAATCCTCGTAAGACTAAGGTTACANCATTCTCCGACCANACTGATCAG
*
ATCAGTCCAGTNATCCTTTTTAGACTTAGGTTGCAACATTCCGTCCGAGCATACTGACCAG
ATGAGACCAGTNATCCTATTGAGACTGAGGTTACATCATTCTCCGACCAAACACTGAACAG
*
ATGAGACCAGTNATCCTNTTNAGACTCAGGTTCCAACATTGCTCCGANCATACTGAGCAG
ATNAGTCCAGTAATCCTTTTAAGACTTAGGTTGAGCATTGGTCCGAGCANACTGACCAG
ATCAGNCCAGTNATCCTATTTAGACTNAGGTTGCATCATTGATCCGATCACACTGANCAG

```

Slika 3: Prikaz grupisanja i sortiranja super-reči

ili različiti ili nedefinirani (označeni sa **N**) u istom bloku. Pvi blok dolazi pre drugog zato što na poziciji označenoj sa zvezdicom se u prvom bloku nalazi **G** nukleotid a u drugom **T** (**G** leksikografski dolazi pre **T**). Slično zvezdica između drugog i trećeg bloka označava zašto drugi dolazi pre trećeg.

2.1 Konstrukcija niza super-reč

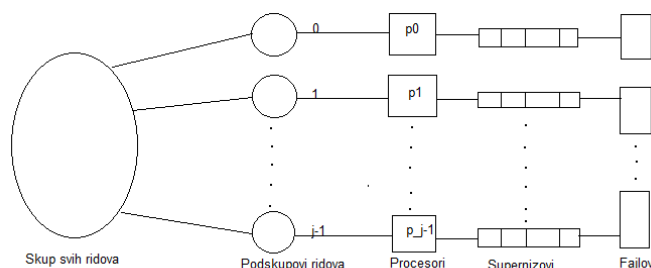
Niz super-reč se konstruiše pomoću bucket sorta u v krugova sortiranja (podsetimo se da se svaka super-reč sastoji od v reči od kojih svaka ima tačno 12 čekiranih pozicija). Najpre se niz inicijalizuje indeksima super-reči u rastućem poretku po vrednostima indeksa. Zatim se iz svake super-reči uzima reč na poziciji p gde p ide od v do 1 (krećemo se sdesna u levo u svakoj super-reči) i niz sortiramo po leksikografskom poretku tih reči. Kako idemo sdesna u levo nakon v koraka bucket sort garantuje da će “niz super-reč” biti sortiran. Za sortiranje se koriste dve pomoćne strukture podataka, tabela pogleda i niz lokacija, koji služe kao kofe (*eng. bucket*) u koje smeštamo reči. Niz lokacija je niz celih brojeva koji je indeksiran od 0 do najvećeg indeksa super-reči. Tabela pogleda se inicijalizuje negativnim brojevima što označava praznu kofu, a ima onoliko polja koliko ima kodova za sve stringove dužine 12. Bucket sortiranje se obavlja tako što se elementi “niza super-reč” sleva na desno čitaju, smeštaju u kofe i zatim se kopiraju iz kofa nazad u “niz super-reč” sdesna u levo. Detaljnije, trenutni element, indeks super-reči, se smešta u kofu iz “niza super-reč” na sledeći način. Indeks super-reči se koristi kako bi se našla reč na poziciji p . Zatim tu reč koja je dužine w transformišemo u string dužine 12 korišćenjem bitovskih operacija kako bi uklonili nukleotide na nečekiranim pozicijama. Znamo da postoji indeks u tabeli pogleda koji predstavlja kod tog stringa dužine 12, pa vrednost tabele pogleda na tom indeksu čuvamo u nizu lokacija na poziciji koja ima isti indeks kao i indeks super-reči, a u tabelu pogleda upisujemo sam indeks super-reči. Kada smo sve elemente “niza super-reč” ubacili u kofe, čitamo ih počevši od najveće kofe i kopiramo ih u “niz super-reč” sdesna ulevo.

2.2 Izračunavanje preklapanja koristeći niz super-reč

Nakon konstrukcije “niza super-reč” delimo ga u sekcije tako da se u svakoj sekciji nalaze identične super-reči. Ove sekcije se obrađuju kako bi se generisala preklapanja između ridova. To radimo na sledeći način. Posmatračemo trenutnu sekciju koja ima s super-reči. Svaka dva rida koja imaju super-reč u ovoj sekciji mogu da se poklapaju. Tako da imamo $(s * (s - 1)) / 2$ potencijalnih preklapanja u ovoj sekciji. Kako bi uprostiti problem mi nećemo računati sva ova preklapanja, već ćemo računati njegov podskup (koji ćemo dobiti za linearno vreme), a ostala preklapanja će moći da budu zaključena iz ovog podskupa. Ovo se može uraditi, tako što sortiramo super-reči u sekciji po dužini njihovih prefiksa u ridu i selektovanjem susednih parova super-reči kako bi izračunali preklapanja. Ako super-reči pripadaju različitim ridovima (podsetimo se da ovo možemo proveriti koristeći funkciju koja od indeksa super-reči računa poziciju super-reči u ridu i indeks rida kome ta super-reč pripada) onda računamo preklapanje između njih. Ukoliko je broj nukleotida koji se razlikuju u preklapanju ograničen odsecanjem, preklapanje čuvamo.

2.3 Ubrzanje paralelizacijom

Veliki skupovi ridova se dele u manje podskupove (obeleženi brojevima od 0 do $j-1$). Svaki rid iz istog podskupa ima super-reč kod koje ostatak pri deljenju koda reči na poziciji p (p je između 1 i v) sa j daje oznaku tog podskupa. Naglasimo da rid može pripadati u više podskupova. Ovi podskupovi se dodeljuju j procesorima kako bi se paralelno izračunala preklapanja. Svaki procesor računa preklapanja u podskupu konstruišući niz super-reč za svaki podskup. Ovim postupkom dobijamo različite fajlove preklapanja za svaki podskup. Ovi fajlovi se potom skeniraju kako bi se uklonila redundantna preklapanja. Na taj način dobijamo skup fajlova neredundantnih preklapanja (slika 4).



Slika 4: Ilustracija paralelnog izvršavanja

2.4 Konstrukcija grafa preklapanja

Fajlovi preklapanja se čitaju redom jedan po jedan i svako preklapanje sa dobrim procentom preklapanja (broj nukleotida koji su se razlikovali u preklapanju je manji od granice) čuvamo u glavnoj memoriji. Graf preklapanja se konstruiše tako da čvorovi budu ridova, a grane između

čvorova su preklapanja između dva reda koju povezuje ta grana. Graf se potom ispituje kako bi se našle dugačke putanje preklapanja neponavljajućih čvorova. Čvor je ponavljajući ako postoje dve dugačke putanje koje se završavaju u čvoru, ali nemaju ponavljanja između njihovih prefiksni putanja. Ovo se izvodi na principu dijekstrinog algoritma. Svaka dugačka putanja neponavljajućih čvorova se predstavlja kao kontiga redova. Stvaranje rasporeda svake kontige se obavlja na jednom procesoru sa velikom količinom glavne memorije. Na kraju ovog koraka se dobijaju fajlovi koji predstavljaju raspored kontiga.

2.5 Generisanje fajla konsenzus sekvenci kontiga i ocene kvaliteta

Fajlovi rasporeda kontiga se paralelno procesiraju. Redovi se raspoređuju kako bi formirali višestruka poravnanja redova. Ova poravnanja se koriste za generisanje konsenzus sekvenci za kontige. Za svaki fajl rasporeda kontiga se generiše fajl konsenzus sekvenci kontiga kao i fajl ocene kvaliteta baza konsenzus kontiga. Na osnovu ovih fajlova se generiše jedan fajl konsenzus sekvenci kontiga i jedan fajl ocene kvaliteta baze konsenzus kontiga.

3 Korišćenje programa PCAP.Solexa

Algoritam opisan u prethodnom poglavlju je implementiran u programu PCAP.Solexa¹. U ovom ovom poglavlju će biti opisan način rada, kao i prateće komande koje je potrebno znati za ovaj program. PCAP.Solexa se pokreće u komandnoj liniji i može se pokrenuti na 64-bitnim Linux/Unix/MacOSX sistemima. Program je dosta memorijski zahtevan, pa za asembliranje genoma gljivice je potrebno 100Gb glavne memorije i 500Gb memorije hard diska. Takođe, potrebno je napomenuti da program nije poželjno pokretati za asembliranje genoma velikih biljaka i životinja.

3.1 Priprema, pokretanje i podaci korišćeni pri radu programa

Nakon skidanja direktorijuma PCAP.Solexa, koji je dostupan na <http://seq.cs.iastate.edu>, može se naći više izvršnih fajlova. Izvršni fajlovi se pokreću po odgovarajućem poretku, koji je zadat u skripti PCAP.Solexa.perl. Ova skripta pre samog pokretanja mora da se modifikuje kako bi stavili do znanja lokaciju PCAP.Solexa direktorijuma i postavili odgovarajuće parametre. U promenljivoj `$CodeDirPath` je potrebno upisati apsolutnu putanju do direktorijuma, kako program mogao da se izvršava. Na primer:

```
$CodeDirPath = "/home/aleksandra/Downloads/pcap.solexa"
```

Ostali parametri su postavljene na podrazumevane vrednosti, opis parametara i vrednosti su dati kao komentari u skripti. U programu je dat jedan primer u poddirektorijumu zvanom test, koji će se koristiti za demonstraciju rada PCAP.Solex programa. On sadrži dva fajla

¹Solexa platforma za sekvenciranje DNK je razvijena od strane britanskih hemičara Shankar Balasubramanian i David Klenerman, a 2006 je komercijalizovana.

`s_4_1_sequence.fastq` i `s_4_2_sequence.fastq` u kojima se nalazi po 7572 rida gde svaki sadrži po 151 nukleotida dobijeni Illumina² sistemom. Poddirektorijum takodje sadrži odgovarajuće `fofn`, `fofn.con` i `fofn.lib` fajlove. Fajl `fofn` sadrži imena fajlova iz kojih se čita `s_4_1_sequence.fastq` i `s_4_2_sequence.fastq`, svaki u posebnom redu. Fajl `fofn.con` ima sledeći sadržaj:

```
s_4_1_sequence.fastq s_4_2_sequence.fastq 100 700 tmp clone
```

Prva dva argumenta označavaju fajlove iz kojih se čitaju ridovi, treći i četvrti donju i gornju granicu (u bp-ovima) koje ćemo unositi, argument `clone` predstavlja ime biblioteke koja će se koristiti, dok je parametar `tmp` čuvar mesta (*eng. placeholder*).

Sadržaj fajla `fofn.lib` je:

```
clone 300 50
```

Drugi i treći argument predstavlja matematičko očekivanje i standardnu devijaciju za biblioteku koja se koristi i koja se nalazi u `clone` argumentu. Pokretanje programa je moguće na dva načina.

Prvi način je komandom:

```
pcap.solexa.perl fofn &
```

čime smo u pozadini pokrenuli rad datoteke `s_4_1_sequence.fastq` i `s_4_2_sequence.fastq`.

Drugi način je pokretanje skripte `cs551_test_script`, koja nas automatski smešta u direktorijum PCAP.Solexa programa i pokreće predhodno navedenu komandu. Pre pokretanja ove skripte je potrebno promeniti apsolutnu putanju do programa PCAP.Solexa.

3.2 Rezultati rada programa

Nakon uspešnog izvršavanja programa, statistike i rezultati asembliranja se nalaze u sledećim fajlovima: `contigs.bases`, `contigsquals`, `reads.placed`, `fofn.pcap.scaffold*.ace`, `readpairs.contigs`, `fofn.con.pcap.results`, `z.n50` i `fofn.pcap.contigs*.snp`.

Fajl `contigs.bases` sadrži sekvence u fasta formatu od 4 kontiga: `Contig0.1`, `Contig1.1`, `Contig2.1` i `Contig3.1`. Brojevi u nazivu kontiga označavaju koja je kontiga u kom skelfodu, recimo `Contig0.1` je prva kontiga u nul-tom skelfodu³ (*eng. scaffold*). Fajl `contigsquals` sadrži Pred ocene kvaliteta sekvenci (*Phred quality score sequence*) u fasta formatu za svaku kontigu (videti dodatak A.1).

Dužine kontiga su date u fajlu `readpairs.contigs`. Sadržaj ovog fajla je sledeći:

```
C0.1 1 23144
C1.1 1 491
C2.1 1 210
C3.1 1 200
```

Prva kolona je skraćeno ime kontiga, druga označava orijentaciju kontiga, a treća dužinu. Dužina kontiga N50 za skup kontiga je definisan kao najveći broj L takav da kontigi dužine barem L imaju ukupnu sumu veću od polovine sume svih dužina kontiga. Broj kontiga N50 za skup kontiga je definisan kao najmanji broj N takav da N-ti najduži kontig ima ukupnu

²Illumina je jedna od poznatijih kompanija za sekvenciranje druge generacije

³Skefold predstavlja uredjeni lanac kontiga orjentisanih jedna posle druge i pomoću kog možemo završiti sekvenciranje genoma.

dužinu veću od polovine sume dužina svih kontiga. Dužina i broj najvećeg kontiga N50 su slično definisani uključivanjem samo kontiga čija je dužina najmanje 1Kb. Statistike kontiga N50 se nalaze u fajlu z.n50. Nakon pokretanja našeg primera u ovom fajlu se nalazi sledeća statistika:

```
Ctg N50 length: 23134, Ctg N50 number: 1
Major ctg N50 length: 23134, Major ctg N50 number: 1
```

Fajl fofn.con.pcap.results sadrži informaciju o tome da li se svaki par uparenih ridova pojavljuje u skefoldu ili kontigi u odgovarajućoj orijentaciji i sa umetnutom veličinom u datom intervalu. Njegov sadržaj izgleda:

```
Ovde ubaci kako izgleda sadržaj fofn.con.pcap.results
nekoliko linija
```

Svaka linija ima barem 6 kolona. Prva i druga kolona predstavljaju indekse dva rida u paru. U ovom primeru indeksi ridova iz fajla s_4_1_sequence.fastq idu od 0 do 7571, a ridovi iz fajla s_4_2_sequence.fastq od 7572 do 15413. Treća i četvrta kolona označavaju donju i gornju granicu intervala ulaza. Kolona pet je čuvar prostora. Šesta kolona je ključna reč koja označava status parova: ključna reč "singlet" znači da barem jedan od ridova iz para nije smešten u asembliju; "redundant" označava da je par isključen jer je bio sličan nekom paru u orijentaciji i rasporedu ridova u asembliranju; "short" označava da je jedan od ridova u paru uključen u kratkom skefoldu ili je pri kraju skefolda. Broj koji se nalazi u šestoj koloni označava da su dva rida uključeni u skalfoldu ili u kontigu u suprotnoj orijentaciji i broj označava razdaljinu između ridova. Ako je razdvojenost unutar intervala ulaza, onda se u sedmoj koloni nalazi opis ridova "zadovoljen u kontigu" (*eng. satisfied in a contig*) ili "zadovoljen u skafoldu" (*eng. satisfied in a scaffold*). Na kraju ovog fajla se nalaze brojeve koji označavaju koliko je parova ridova u svakoj kategoriji pročitano.

U fajlovima **fofn.pcap.contigs*.snp** se nalaze izveštaji o kandidatima za polimorfizam pojediničnih nukleotida (*single nucleotide polymorphisms (SNP)*)⁴. Broj tih fajlova je vrednost V za parametar **\$NoCutJobs** koja se nalazi u skripti pcap.solexa.perl, gde su fajlovi indeksirani od 0 do $V-1$. Fajl **fofn.pcap.cotignsN.snp** sadrži izveštaj SNP-ova u skefoldima N , $N + V$, $N+2 * V$ itd.

```
Prikaži par redova fofn.pcap.cotignsN.snp
i opiši šta znači koja kolona.
```

Izveštaj o SNP-u se nalazi u linijama koje počinju sa "SP". Kolone u ovim linijama imaju sledeća značenja:

- kolona 3: pozicija SNP-a u kontigi
- kolona 4: ime kontige
- kolona 2: dubinu pokrivenosti SNP-a, odnosno broj ridova obeleženih sa BS u ostatku sekcije.

Kolone u "BS" linijama označavaju višestruka poravnanja ridova na SNP pozicijama u kontigi. Svaki rid u poravnanju je prikazan u ovoj liniji sa sledećim značenjima:

- kolona 2: baza rida u SNP poziciji
- kolona 3: ocena kvaliteta baze
- kolona 4: dužina rida

⁴SNP je varijacija jednog nukleotida koji se javlja na određenoj poziciji u genomu, gde je svaka varijacija prisutna do nekog stepena unutar populacije.

- kolona 5: ime rida

U fajlu `reads.placed` se nalazi izveštaj o poziciji svakog rida u asembliranju. Fajl je organizovan po skofildima (koji se nazivaju i superkontige). Za svaki rid u kontigi u skefoldu kolone imaju sledeća značenja:

- kolona 2: ime rida
- kolona 4: dužina rida
- kolona 5: bit koji označava orijentaciju (0 je data orijentacija) rida u kontigi
- kolona 6: ime kontige
- kolona 7: ime skefolda
- kolona 8: pozicija rida u kontigi
- kolona 9: pozicija rida u skefoldu
- ostale kolone su čuvari mesta.

3.3 Potencijalni problemi

U slučaju da u promenljiva `$CodeDirPath` u skripti `pcap.solexa.perl` nije definisana puna putanja do `pcap.solexa` direktorijuma, sledeća greška će biti prikazana prilikom pokretanja programa:

```
The definition ("'/home/xqhuang/551/pcap.solexa'") for the variable $CodeDirPath in this Perl script is incorrect.
```

U slučaju da promenljive `$NoPcapJobs` i `$NoCutJobs` u skripti `pcap.solexa.perl` su postavljene na vrednost koja je veća od broja skafolda, pojavice se greške jer će broj izlaznih skofild fajlova biti veći od samog broja skofilda. Konačno, u slučaju da ulazni fajlovi sadrže grešku, `PCAP.Solexa` će izbaciti veoma čudnu grešku. Zbog toga je bitno da se pre pokretanja programa provere ulazni fajlovi, da bi bili sigurni da je sve u redu.

4 Zaključak

U radu smo opisali i demonstrirali način rada programa `PCAP`. Pomenuli smo nekoliko karakteristika koje se odnose na efikasnost i problem tačnosti asembliranja u programu. Višestruki procesori se koriste kako bi se obavila izračunavanja koja troše dosta vremena. Osetljiviji metodi su korišćeni kako bi se izbegla nedostajuća preklapanja nastala usled grešaka sekvenciranja. Videli smo da je stvaranje konsenzus sekvenci kontiga bazirano na poravnanju ridova u kontigama u kojima se vrednost kvaliteta baze i informacija o pokrivenosti koriste da odrede svaku bazu konsenzusa. Ovaj program je pokazao dosta dobre rezultate, a neka od testiranja su obavljena na skupu podataka celog genoma miša koji je imao 30 miliona ridova, kao i na 20 skupova podataka od po 1.7 miliona ridova ljudskog hromozoma.

Literatura

- [1] Jonathan M. Keith. *Bioinformatics Volume I: Data, Sequence Analysis and Evolution*.
- [2] Philip Compeau, Pavel Pevyner *Bioinformatics Algorithms* (2nd Edition, Vol. 1)

- [3] PCAP: A Whole-Genome Assembly Program,
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC403719/>

A Dodatak

A.1 Ocena kvaliteta sekvenci Phred

Sekvenciranje ridova sa lošim kvalitetom utiče na tačnost asembliranja. Zbog toga filtriranje podataka dobrog kvaliteta iz svih podataka sekvenciranja je potrebno. Jedna od najčešćih ocena koja se koristi pri ocenjivanju kvaliteta je Phred. Ona definiše vrednost kvaliteta q na sledeći način:

$$q = -10 * \log_{10}(p)$$

gde p predstavlja procenjenu verovatnoću greške. Tako na primer verovatnoći da je netačnost 1/1000 je dodeljena ocena Phred 30 (slika5).

Phred quality score	Error rate	Accuracy of base call
10	1 in 10	90 %
20	1 in 100	99 %
30	1 in 1000	99.9 %
40	1 in 10,000	99.99 %
50	1 in 100,000	99.999 %

Slika 5: Ocena verovatnoće phred ocene