

Penggunaan Paralelisme dan Terdistribusi untuk Pemrosesan *Phylogenetic Tree* pada *Sequence Biologi* dengan Menggunakan *MapReduce*

Renaning Karutami Susilo¹, Setyorini², Siti Amatullah Karimah³

^{1,2,3}Fakultas Informatika, Universitas Telkom, Bandung
¹renaaks@students.telkomuniversity.ac.id, ²setyorini@telkomuniversity.ac.id,
³karimahsiti@telkomuniversity.ac.id

Abstrak

Multiple Sequence Alignment (MSA) merupakan proses penting dalam analisis *sequence biologi* dengan melakukan perbandingan pada sejumlah *sequence biologi*. Pada beberapa algoritma MSA (seperti pada CLUSTALW, misalnya), pembentukan *phylogenetic tree* sebagai *guidence* dalam proses *alignment* memiliki peran penting dalam menentukan akurasi hasil akhir *alignment*. Dari keseluruhan proses MSA, pembentukan *phylogenetic tree* memerlukan waktu komputasi yang meningkat seiring dengan peningkatan jumlah *sequence*. Komputasi score similaritas untuk semua kombinasi pasangan *sequence* yang dilaksanakan secara berurutan menjadi masalah pada waktu komputasi. Tugas Akhir ini meneliti potensi efisiensi komputasi *phylogenetic tree* secara paralel dan terdistribusi pada lingkungan Hadoop menggunakan *MapReduce*. Hasil penelitian menunjukkan *phylogenetic tree* dapat dibangun dengan menggunakan komputasi *MapReduce*.

Kata kunci : Biology sequence, Multiple Sequence Alignment, Phylogenetic Tree, Hadoop, *MapReduce*.

Abstract

Multiple Sequence Alignment (MSA) is an important process in the analysis of biological sequences by making comparisons of a number of biological sequences. In some MSA algorithms (such as CLUSTALW, for example), the formation of phylogenetic trees as a guideline in the alignment process has an important role in determining the accuracy of the final alignment results. From the whole MSA process, the formation of phylogenetic tree computational time increased as the number of sequences increases. Computational score similarity for all combinations of sequence pairs carried out sequentially becomes a problem at computational time. This Final Project examines the potential efficiency of computational phylogenetic trees in parallel and is distributed to the Hadoop environment using *MapReduce*. The results showed that phylogenetic tree can be generated using *MapReduce* computation.

Keywords: Biology sequence, Multiple Sequence Alignment, Phylogenetic Tree, Hadoop, *MapReduce*.

1. Pendahuluan

Latar Belakang

Asam dioksibirinukleat atau yang secara umum dikenal dengan DNA adalah molekul biologis yang menyimpan berbagai informasi genetika sebuah organisme. DNA merupakan asam nukleat, salah satu molekul yang hakiki bagi makhluk hidup yang ada. DNA umumnya memiliki dua untai polinukleotida yang masing-masing terdiri dari nukleotida dengan salah satu jenis basa nitrogen (guanina, adenine, timina atau sitosina) yang saling terikat menjadi satu rantai [1]. Dari rantai ini didapatkan kombinasi *string* yang berbeda-beda. Kombinasi *string* ini menyimpan banyak informasi mengenai organisme tersebut, contohnya seperti fungsi, keturunan, struktur dan lainnya [1]. Dalam bioinformatika, analisis kemiripan antar *sequence* dilakukan dengan metode *sequence alignment*. Informasi kemiripan pola *sequence* membawa kemiripan pada fungsi, keturunan, struktur dan informasi genetik [1]. Algoritma untuk melakukan *sequence alignment* diantaranya yaitu algoritma Needleman-Wunsch yang menghasilkan *global alignment* dan algoritma Smith-Waterman yang menghasilkan *local alignment*. *Global alignment* mencari nilai kecocokan tertinggi dari ujung awal DNA hingga ujung akhir DNA, sehingga hasil *alignment* akan sama panjang, sementara *local alignment* hanya mencari hasil kecocokan tertinggi berdasarkan bagian DNA tersebut [2].

Untuk kebutuhan membandingkan *sequence* dalam lebih dari 2 menggunakan metode *multiple sequence alignment (MSA)*. Salah satu contohnya Clustal yang menggunakan algoritma Needleman-Wunsch pada proses di dalamnya [3], ada juga T-COFFEE (*Tree Based Of Consistency Objective Function For Evaluation Alignment*) proses didalamnya menggunakan gabungan dari algoritma Smith-Waterman dan juga Needleman-Wunsch. Hasil T-COFFEE lebih akurat dibandingkan dengan Clustal karena menggunakan *local* dan *global alignment* [4]. Pada MSA dibangun sebuah *phylogenetic tree* (pohon evolusi) yang menunjukkan kedekatan tiap sekuens DNA yang ada dan merepresentasikan hubungan evolusi dalam organisme tersebut. Namun, algoritma yang digunakan untuk MSA membutuhkan banyak waktu komputasi yang lama [2]. Untuk mengatasi masalah ini, eksekusi algoritma

MSA dapat dilakukan secara paralel terdistribusi. Hadoop merupakan sebuah *framework* yang membantu dalam melakukan proses secara terdistribusi dan paralel [5]. Di dalam *framework* Hadoop terdapat sebuah metode bernama *MapReduce* dimana metode ini dapat melakukan proses paralel dari data-data berukuran besar [5]. Algoritma MSA dapat memanfaatkan *framework* Hadoop dan juga metode *MapReduce* untuk mempercepat waktu komputasi walaupun data yang diproses berukuran besar [6]. Dalam melakukan melakukan proses terhadap DNA dapat digunakan *MapReduce* yang bertugas untuk melakukan *sequence alignment* pada tiap baris DNA dari hasil tersebut kemudian dapat digunakan untuk melihat kedekatan hubungan antar baris pada DNA.

Topik dan Batasannya

Berdasarkan hal-hal yang sebelumnya telah dibahas di latar belakang, berikut ini merupakan rumusan-rumusan masalah yang akan dibahas dalam penelitian ini:

1. Bagaimana membangun *phylogenetic tree* menggunakan model komputasi *MapReduce*?
2. Proses mana yang dapat diparalelkan dalam pembangunan *phylogenetic tree*?

Adapun batasan-batasan yang ditetapkan dalam menjawab rumusan masalah diatas, yaitu:

1. Menggunakan *framework* Hadoop dengan *MapReduce* dan Hadoop File System (HDFS) .
2. Metode yang digunakan untuk membangun *phylogenetic tree* adalah NJ (Neighbor Joining).

Tujuan

Tujuan dibuatnya tugas akhir ini adalah untuk membangun *phylogenetic tree* dengan menggunakan metode NJ (*Neighbor Joining*) dengan Hadoop *MapReduce* sebagai model komputasi.

Organisasi Tulisan

Pada jurnal ini, terdiri dari lima bab utama, yang pertama merupakan Pendahuluan. Bab kedua merupakan Studi Terkait. Bab ini berisi teori-teori penunjang penelitian yang dilakukan, yang dibahas disini meliputi Bioinformatika, Algoritma-algoritma yang digunakan dalam pemrosesan sekuens biologi, Penjelasan mengenai *framework* Hadoop dan juga model komputasi *MapReduce*. Selanjutnya bab ketiga yaitu Sistem yang Dibangun, pada bab ini dijelaskan mengenai rancangan sistem yang akan dihasilkan. Bab keempat yaitu Evaluasi. Pada bab ini terdiri dari hasil rancangan sistem yang telah dibangun yang dilakukan sesuai rancangan sistem pada bab tiga. Terakhir yaitu bab kelima, berisi tentang kesimpulan dari hasil pengujian sistem yang dibangun dan juga saran untuk penelitian yang akan dilakukan selanjutnya.

2. Studi Terkait

Sequence Alignment

Bioinformatika adalah ilmu yang mempelajari mengenai penerapan teknik komputasional dalam informasi di bidang biologi. [7] Pada bidang ini digunakan penerapan metode matematika, statistika dan informatika dalam memecahkan berbagai permasalahan biologi. Salah satu masalah populer pada bidang bioinformatika adalah mengenai deretan DNA atau protein lainnya. [7]

Multiple Sequence Alignment atau yang biasa disebut MSA adalah cara menyusun urutan baris DNA untuk mengidentifikasi suatu kesamaan sifat, fungsi dan informasi genetik lainnya dengan melibatkan lebih dari dua buah baris DNA, sehingga disebut *multiple*. MSA ini memiliki berbagai macam algoritma, contohnya adalah algoritma Needleman-Wunsch dan algoritma Smith- Waterman. [2]

1. Algoritma Needleman-Wunsch

Algoritma ini digunakan untuk melakukan *global alignment* dimana panjang sekuens seluruhnya terlibat dalam perhitungan yang dilakukan untuk menemukan *alignment* yang optimal dari sepasang DNA ini. [2]

```

ACTACTAGATTACTTACGGATCAGGTTACTTTAGAGGCTTGCAACCA
||||| ||||| |||  ||||| |||  ||||| ||||| ||||| ||||| |||||
ACTACTAGATT - - - -ACGGATC- -GTACTTTAGAGGCTAGCAACCA
    
```

Gambar 1 Contoh Global Alignment

2. Algoritma Smith-Waterman

Algoritma ini digunakan untuk melakukan *local alignment* dimana hanya sebagian dari panjang sekuens yang digunakan dalam perhitungan untuk menemukan *alignment* yang optimal. [2]

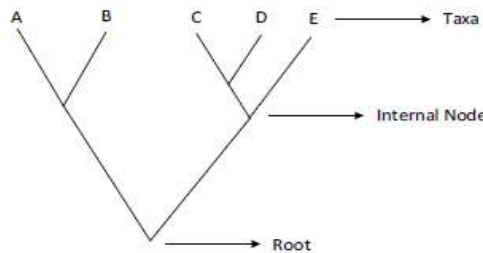
```

ACTACTAGATTACTTACGGATCAGGACTTTAGAGGCTTGCAACCA
||||| ||||| ||||| ||||| |||||
TACTCACGGATGAGGACTTTAGAGGC
    
```

Gambar 2 Contoh Local Alignment

Phylogenetic Tree

Hasil dari algoritma diatas merupakan skor yang merepresentasikan kedekatan dari dua *sequence* DNA dan juga bentuk *sequence alignment* yang baru, skor ini dapat digunakan untuk membangun sebuah pohon kedekatan atau yang biasa disebut *phylogenetic tree*. *Phylogenetic Tree* atau Pohon Filogenetika adalah diagram percabangan (*tree*) yang menunjukkan hubungan evolus atau kedekatan antar berbagai spesies makhluk hidup. [8] Pohon Filogenetika ini dapat digunakan untuk membuat sistematika biologi, contohnya adalah pohon kehidupan. Selain itu, pohon ini dapat berperan untuk mencari fungsi/sifat suatu gen atau protein.



Gambar 3 *Phylogenetic Tree*

Berikut bagian-bagian yang ada pada *phylogenetic tree*:

1. Root: Merupakan akar awal terjadinya percabangan pada sebuah tree.
2. Taksa: Bagian ujung pada sebuah cabang, merepresentasikan spesies atau sekuens
3. Internal Node: Titik yang menghubungkan dua cabang, node merepresentasikan keturunan. [9]

Dalam membangun *phylogenetic tree* dapat menggunakan berbagai metode. Metode yang banyak digunakan adalah UPGMA (Unweighted Pair Group Method using Arithmetic Average) dan NJ (Neighbor Joining)

1. UPGMA (Unweighted Pair Group Method using Arithmetic Average)

Menggunakan sekuensial clustering method. Menggabungkan 2 taksa dengan nilai *pairwise distance* terkecil pada *matrix distance*. Kemudian dihitung lagi dengan menempatkan node ditengah antara dua taksa itu, kemudian dihitung lagi nilainya sehingga terbentuk matriks yang telah tereduksi. [9]

Misalkan ada 4 buah sekuens biologi yaitu A, B, C dan D, dengan menggunakan rumus kombinasi:

$$C(4,2) = \frac{4!}{(2!(4-2)!)} = \frac{24}{4} = 6$$

maka terdapat 6 pasangan sekuens. Berikut adalah himpunan pasangan sekuens dengan masing-masing skornya $\{(AB,3),(AC,3),(AD,6),(BC,7),(BD,8),(CD,9)\}$.

a. Membangun matriks

Ukuran matriks yang dibangun sama dengan jumlah sekuens yang dimiliki, contohnya disini ada 4 sekuens maka matriks yang dibangun berukuran 4x4. Matriks ini kemudian diisi dengan skor tiap pasangan sekuens berdasarkan himpunan skor.

	A	B	C	D
A	0	5	3	6
B	5	0	7	8
C	3	7	0	9
D	6	8	9	0

Gambar 4 Matriks skor kedekatan

b. Mencari *alignment* terdekat dalam matriks

Dari matriks yang telah dibangun, cari dua *alignment* yang terdekat, dilihat dari skornya yg terkecil, kemudian dua *alignment* dibentuk menjadi satu node karena cabang *equidistant* maka skor AC dibagi dua.

$$\frac{AC}{2} = \frac{3}{2} = 1.5$$



Gambar 5 Nilai terkecil pada matriks bobot dan node yang dibentuk

c. Reduksi matriks dan hitung skor terbaru

Setelah mendapatkan skor terkecil maka matriks direduksi menjadi seperti pada gambar 7. Kemudian nilai jarak untuk B dan D menuju AC dihitung kembali dengan rumus:

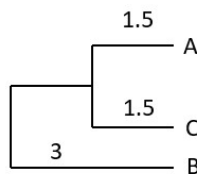
$$ACB = \frac{AB + BC}{2} = \frac{5 + 7}{2} = \frac{12}{2} = 6$$

$$ACD = \frac{AD + CD}{2} = \frac{6 + 9}{2} = \frac{15}{2} = 7.5$$

	AC	B	D
AC	0	6	7.5
B	6	0	8
D	7.5	8	0

Gambar 6 Matriks yang telah direduksi dan skor terbaru

Setelah didapatkan matriks terbaru seperti pada gambar 8, lakukan hal yang sama seperti pada langkah b. Skor terkecil adalah 6 yaitu dari AC ke B maka dibentuk cabang baru dari AC ke B seperti pada gambar 9. Karena *equidistant* maka skor 6 dibagi 2, $\frac{6}{2} = 3$.

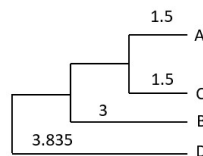


Gambar 7 Cabang B yang telah ditambahn ke node AC

Ulang langkah b dan c hingga seluruh *alignment* terbentuk cabangnya. Disini tersisa jarak ACB ke D maka dihitung dengan rumus

$$\frac{AD + CD + BD}{3} = \frac{8 + 6 + 9}{3} = 7.67$$

karena *equidistant* maka 7.67 dibagi 2, $\frac{7.67}{2} = 3.835$.



Gambar 8 Bentuk akhir tree yang terbentuk

2. NJ (Neighbor Joining)

Mirip dengan UPGMA dalam membangun tree menggunakan *stepwise reduced matrix*. Namun pada NJ nilainya tidak diasumsikan ekuivalen, namun ada *evolutionary rate* pada tiap sekuens dengan menggunakan *conversion step*. Pada *conversion step* digunakan nilai r dan juga nilai *transformed r*. [9] rumusnya seperti dibawah ini:

$$d'_{AB} = d_{AB} - 1/2 \times (r_A + r_B)$$

d'_{AB} adalah *converted distance* antar sekuens A dan B

d_{AB} adalah nilai *actual evolutionary distance* antar sekuens A dan B

$r_A + r_B$ adalah jumlah jarak A atau B pada taksa lainnya, nilai r didapat dari rumus:

$$r_i = \sum d_{r_{ij}}$$

Dimana i dan j adalah dua taksa yang berbeda, nilai *transformed r* dibutuhkan untuk menentukan jarak dari individual taksa ke node terdekat

$$r'_i = \frac{r_i}{n} - 2$$

Dimana ada n buah taksa, misalnya A dan B berasal dari node U, maka jarak dari A ke u ditentukan dengan rumus:

$$d_{AU} = \frac{[d_{AB} + (r'_A - r'_B)]}{2}$$

Misalkan ada 3 buah sekuens DNA yaitu A, B dan C, dengan menggunakan rumus kombinasi:

$$C(3,2) = \frac{3!}{(2! (3 - 2)!)} = \frac{6}{2} = 3$$

maka terdapat 3 pasangan sekuens. Berikut adalah himpunan pasangan sekuens dengan masing-masing skornya $\{(AB,3),(AC,3),(BC,7)\}$.

a. Membangun Matriks

Ukuran matriks yang dibangun sama dengan jumlah sekuens yang dimiliki, contohnya disini ada 3 sekuens maka matriks yang dibangun berukuran 3x3. Matriks ini kemudian diisi dengan skor tiap pasangan sekuens berdasarkan himpunan skor.

	A	B	C
A	0	5	3
B	5	0	7
C	3	7	0

Gambar 9 Matriks skor Kedekatan

b. Menghitung nilai *r-value* dan *r'-value*

Hitung nilai *r-value* dan *r'-value* dihitung dengan menggunakan rumus dibawah ini:

$$r_A = AB + AC = 5 + 3 = 8$$

$$r'_A = \frac{r_A}{3 - 2} = \frac{8}{1} = 8$$

$$r_B = BA + BC = 5 + 7 = 12$$

$$r'_B = 13$$

$$r_C = CA + CB = 3 + 7 = 10$$

$$r'_C = 10$$

Nilai *r'-value* sama dengan nilai *r-value* karena pembagiannya adalah 1.

c. Menghitung nilai *converted distance*

Hitung nilai *converted distance* dengan menggunakan rumus dibawah:

$$d'_{AB} = d_{AB} - \frac{1}{2} * (r_A + r_B) = 5 - \frac{1}{2} * (8 + 12) = -5$$

$$d'_{AC} = d_{AC} - \frac{1}{2} * (r_A + r_C) = 3 - \frac{1}{2} * (8 + 10) = -6$$

$$d'_{BC} = d_{BC} - \frac{1}{2} * (r_B + r_C) = 7 - \frac{1}{2} * (12 + 10) = -4$$

Ganti nilai matriks menggunakan nilai yang didapatkan pada langkah ini, sehingga matriks menjadi

	A	B	C
A	0	-5	-6
B	-5	0	-4
C	-6	-4	0

Gambar 10 Matriks dengan nilai *converted distance*

d. Pilih nilai terkecil

Dari matriks pada gambar, pilih nilai terkecil, didapatkan -6 dari A ke C yang kemudian akan digabungkan menjadi satu *node tree* yang diberi nama U.

	A	B	C
A	0	-5	-6
B	-5	0	-4
C	-6	-4	0

Gambar 11 Nilai terkecil pada matriks bobot dan node yang dibentuk

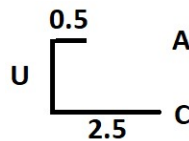
e. Hitung jarak tiap cabang

Jika pada metode UPGMA jarak cabang dianggap sama, maka pada metode NJ jarak cabang dihitung berbeda, jarak dihitung masing-masing menggunakan rumus

$$d_{AU} = \frac{[d_{AC} + (r'_A - r'_C)]}{2} = \frac{[3 + (8 - 10)]}{2} = 0.5$$

$$d_{CU} = \frac{[d_{AC} + (r'_C - r'_A)]}{2} = \frac{[3 + (10 - 8)]}{2} = 2.5$$

Dari hasil perhitungan didapatkan AU = 0.5 dan CU = 2.5, maka cabang digambarkan sebagai berikut



Gambar 12 Tree yang terbentuk

Panjang cabang berbeda sesuai dengan nilai yang didapatkan.

f. Menghitung jarak baru

Karena A dan C sudah menjadi satu node, maka nilainya dengan B harus diperbarui dengan menggunakan rumus

$$d_{BAC} = \frac{[(d_{AB} - d_{UA}) + (d_{BC} - d_{UC})]}{2} = \frac{[(5 - 0.5) + (7 - 2.5)]}{2} = 4.5$$

Hadoop & MapReduce

Pada sistem ini akan digunakan *framework* Hadoop dan juga metode *MapReduce* dalam melakukan pemrosesan sekuens biologi. Hadoop merupakan sebuah *framework* yang mensupport aplikasi atau pengolahan yang melibatkan big data [10], Hadoop melakukan komputasi dengan distribusi ke node. Hadoop dapat dijalankan secara *single-node* dan *multi-node*. Pada *single-node*, seluruh kegiatan dilakukan hanya menggunakan satu perangkat saja sedangkan *multi-node* dilakukan pada banyak perangkat yang ada pada kluster. Kluster ini kumpulan dari perangkat yang saling terhubung dan bekerja dalam Hadoop atau biasanya disebut sebagai Hadoop Cluster. Dalam Hadoop Cluster *multi-node* biasanya ada perangkat yang berperan sebagai kontrol utama dalam proses (*Master*) dan perangkat yang berperan menjalankan proses yang diberikan (*Slave*). Pada Hadoop perhitungan atau proses komputasi itu bisa dibagi ke perangkat lainnya tanpa saling tindih sehingga hasil menjadi cepat semakin banyak perangkat yang terlibat dalam proses maka akan semakin cepat proses tersebut dilakukan.

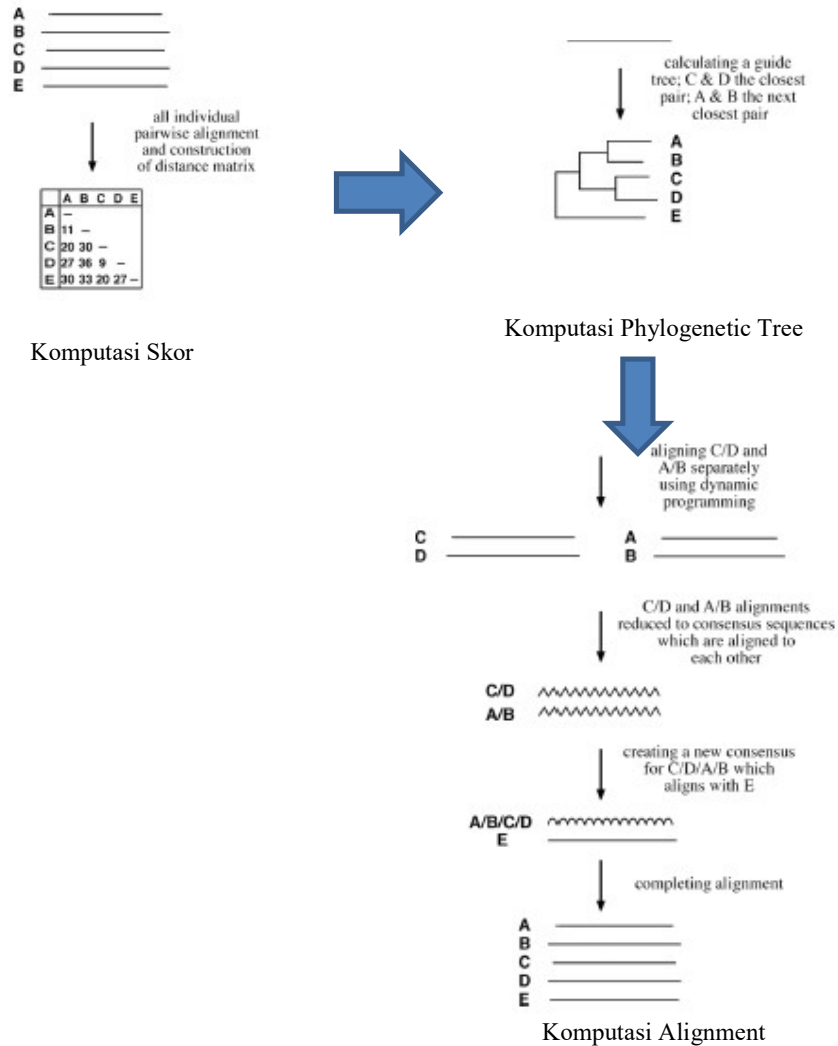
Hadoop memiliki *file system* sendiri yang disebut *Hadoop Distributed File System* (HDFS). HDFS menyimpan data yang digunakan dalam proses yang berjalan di Hadoop Cluster. Dengan adanya HDFS Sistem tersebut dapat membagi beban penyimpanan ke berbagai perangkat, walaupun menggunakan perangkat yang berbeda, namun jika salah satu perangkat mati tidak akan berdampak apapun pada perangkat lain, data akan masih tersimpan dan proses masih bisa terus dijalankan. [10] Ini merupakan salah satu contoh *failure control* yang dimiliki Hadoop. HDFS pada Hadoop *multi-node* juga memiliki kontrol utama seperti *master-slave*, namun dalam HDFS yang berperan sebagai *master* disebut *name node* dan yang berperan sebagai *slave* disebut *data node*.

MapReduce merupakan sebuah metode komputasi yang tersedia dalam *framework* Hadoop. *MapReduce* ini memiliki *job tracker* dan *task tracker* yang memiliki peran mirip dengan *master-slave* pada Hadoop. *MapReduce* pada sistem terdistribusi memiliki dua proses utama yaitu *map* dan *reduce*. *Map* berfungsi untuk membagi sebuah proses menjadi beberapa bagian yang dapat diproses secara paralel dengan jumlah yang banyak juga dengan ukuran data yang besar, dimana tiap node menerapkan fungsi *map* ke lokal data dan menuliskan hasil prosesnya ke penyimpanan sementara. [6] *Reducer* berfungsi sebagai penggabung hasil proses yang dilakukan secara terpisah oleh banyak *map*. Fungsi *map* dan *reduce* didefinisikan dengan menggunakan sepasang *key* dan *value*. [7] *Map* mengambil sepasang data (*key, value*) pada satu data domain dan mengembalikan sebagai *list of pair* pada domain yang berbeda. [6] Semua yang dilakukan di dalam Hadoop *MapReduce* terjadi secara otomatis, user hanya perlu menjalankan program java dalam bentuk file jar.

3. Sistem yang Dibangun

3.1 Analisis Independen Proses pada Pembentukan *Phylogenetic Tree*

Untuk kebutuhan paralel dan distribusi diperlukan identifikasi proses independen pada komputasi MSA. Proses independen adalah unit proses yang tidak saling bergantung satu dengan yang lainnya. Proses-proses ini yang berpotensi untuk dijalankan secara paralel dan terdistribusi. Pada penelitian ini mengambil contoh algoritma MSA Clustal W. Berikut adalah proses MSA dari CLUSTAL W

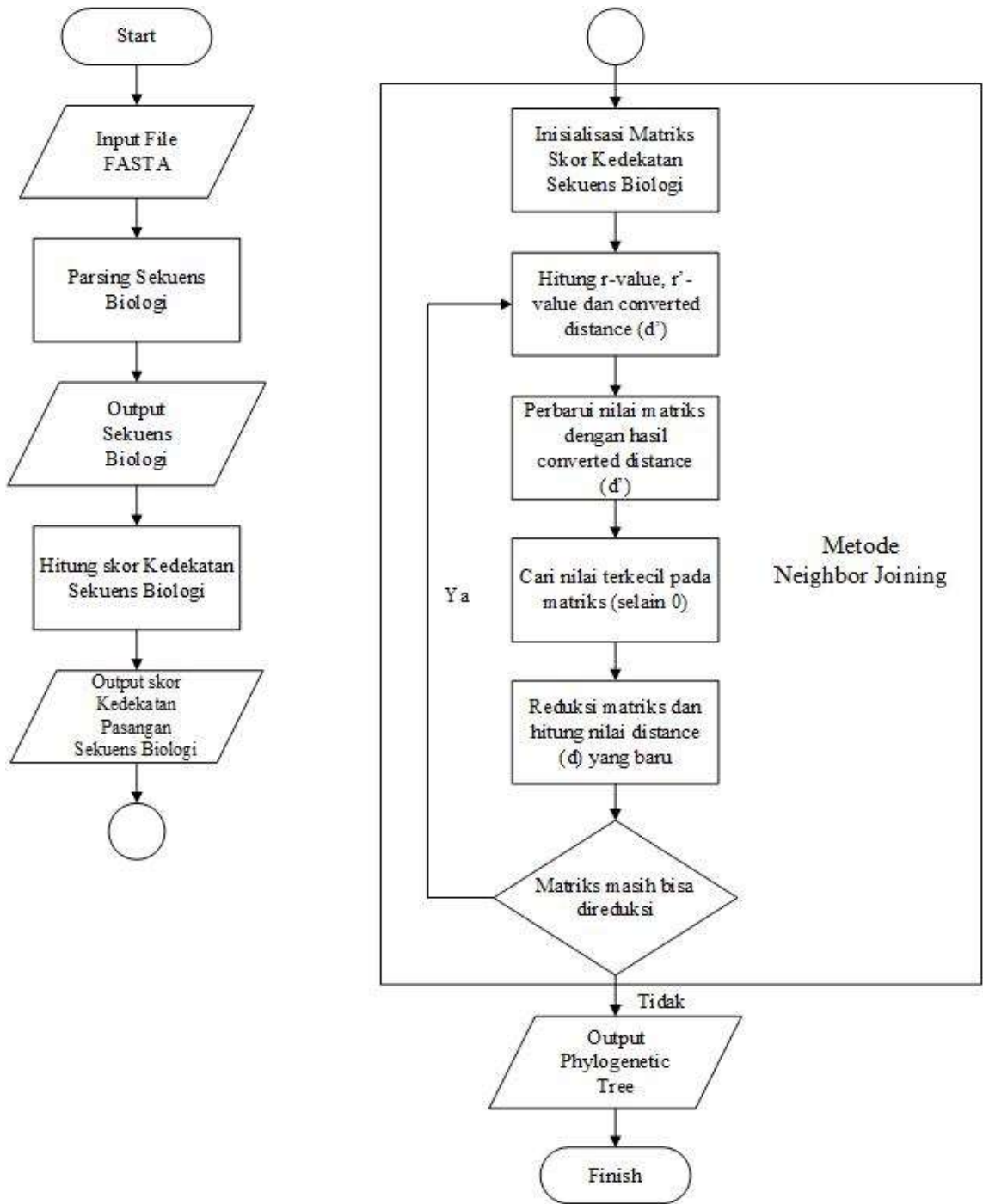


Gambar 13 Alur proses MSA (CLUSTAL W)

Dari gambar 13 komputasi skor memiliki potensi melibatkan banyak proses independen, karena komputasi skor dilakukan untuk semua kombinasi pasangan *sequence* yang terbentuk dari sejumlah input *sequence*. Masing-masing komputasi skor dapat dipandang sebagai unit proses yang independen.

3.2 Rancangan Algoritma Phylogenetic Tree

Tahapan proses pembangunan *phylogenetic tree* dapat digambarkan sebagai berikut:



Gambar 14 Alur proses pembentukan *phylogenetic tree*

Waktu komputasi yang digunakan dalam membangun *phylogenetic tree*, waktu yang ditunjukkan adalah 2064 ms. Waktu komputasi ini didapatkan dari menghitung selisih waktu saat program ini dijalankan hingga program menghasilkan *phylogenetic tree*.

```

33234702):90.34476649204146,AD1753_PANIK/227
76604):72.94718097150326):9.084725519472897,
-337:525.9024551746458,(((Q626WE_CAEBR/159-3
49:435.27214617768595):16.64221722146739):43
,Q61FX7_CAEBR/161-362:473.1820020488664):24.
89_CAMDR/7-52:609.3358587537493):18.53671722
9292035):36.242655645940715,NH174_CAEL/233-
):0.7854039710946381,(((Q9GTC8_CAEL/193-3
1-424:457.43388500548247):37.65374940814394)
8,NHR18_CAEL/183-216:598.6626466320407):25.
993929227193,Q0GGP6_GAMHO/6-86:614.100607077
Tree Construction: 2064 ms.
    
```

Gambar 17 Waktu komputasi yang digunakan IDE Java

Namun, tiap menjalankan program waktu yang dihasilkan berbeda-beda, berikut waktu yang dihasilkan dalam 4 kali menjalankan program.

Tabel 1 Rata-rata waktu komputasi yang digunakan IDE Java

Run Program ke-	Waktu yang digunakan (ms)
1	2064
2	2649
3	3013
4	2770
Rata-rata	2624

Rata-rata waktu yang digunakan dalam membentuk tree adalah 2624 ms.

- Hasil Parsing file FASTA dan pembentukan *phylogenetic tree* dengan menggunakan *MapReduce* yang dijalankan pada Hadoop multinode, dengan jumlah node sebanyak dua buah. Berikut *phylogenetic tree* dengan format newick

```

1 1
((((((A2D504_ATEGE/1-46:555.3480174177785,(((O18972_BOVIN/55-125:447.9152227227272,PRGR_MACEU/107-178:448.0847772272728):47.87301788330078,((Q
410102_9PABR/125-230:349.4947017840485,Q4WIT7_BOOLU/118-240:347.5052982159915):22.125866117931537,Q2PR04_SMOBO/114-214:375.37413388206846):118.6
269021165922):20.46525949360563,AL1YMS_SMOLE/1-45:536.7047405635949):11.5736386170249,AD2SQ1_LAGLA/1-46:550.8013613382975):3.151982592221535
3):42.6733836954027,((((((Q63879_SMRRI/1-161:270.1706627155172,(((A7X809_CERAE/723-912:148.0_A7XW20_PITIR/725-914:148.0):0.0_A7X809_PONVY/7
25-914:148.0):0.0_A7XSD2_COLGU/724-913:148.0):0.0_A7XW16_CEBAE/725-914:148.0):1.0172413793103487,A7X8C2_HYLLA/722-911:149.58275862068965):17.84
47265625,(A7X8D4_TRAOB/723-912:146.5_A7XW25_ATEPA/725-914:146.5):19.6552734375):50.3293728448279):29.069122779187822,PRGR_RANDY/501-690:248.43
08772208122):90.54470687984497,O19008_HORSE/1-118:393.93029312015506):58.77247822392094,((((((MCR_TUGGB/767-956:176.23314606741573,MCR_SAIBC/77
2-961:156.7689393258427):60.78983369098712,MCR_PIG/1-154:275.2101663090129):29.522648358585858,(A2CE98_DANRE/162-351:168.64092664092664,Q5WF01
_ONCMY/528-1017:175.35907335907336):87.4735164141415):86.22713341346154,Q90DM1_PEMA/228-417:334.39786658653844):26.158464862454323,(((A0P3P7
_FORNO/93-281:168.36046511627907,ATLAC6_SICCH/99-287:163.6393488372093):20.591412401574804,(Q6BRQ3_ONCMY/460-648:157.54131615384615,Q2W9E_PU
AFE/144-332:157.45865384615385):9.85749027237354,Q5I124_SALTR/539-727:170.64250972762648):17.533587598425156):34.73556022970085,Q1L691_DANRE/12
6-310:224.32693977029913):44.78402363885309,(Q2E831_XENTRY/569-757:242.0515625,(Q33Y1_SMIMA/34-218:179.89462594117645,(((GCR_CALJA/568-756:139.5
_GCR_SAGOR/568-756:139.5):0.0_GCR_AOTNA/568-756:139.5):17.123134328358205,GCR_MOUSE/574-762:158.37686567164178):20.105147058823547):64.4484375
000001):24.208163861146907):103.66184763754568):20.556580947617352,((((((A4LAN9_SPRIM/708-897:132.5_ANDR_BULFC/674-863:132.5):23.517867647058
84,Q6E9Y8_COEJA/133-322:158.98621323529412):47.68876518218623,B70048_XENLA/580-769:199.81123481781378):71.0782953910615,Q5EVO2_SQOUC/568-757:2
68.92637046009395):19.60811335421683,Q9DD34_BULTR/354-539:301.016948445793):24.55335031147876,(Q600F9_DANAP/474-659:196.13191699604744,Q9W6P4
_HABBU/481-666:159.860808300395256):125.60289568952123):72.26861436488265):18.98759119829058):35.30668584160182,Q28854_MACMU/4-108:459.43501337
714815):76.61737120778938,Q6LWJ2_XENLA/1-70:573.7869256672107):34.1152820819441):15.387733039358864,Q988T1_CHICE/15-61:624.8263325123991):11.940
56906328587,((((((((((Q5NOR2_PIG/1-87:456.7798563639323,Q1RFA3_BUBBU/2-65:477.2201436360677):42.23454266283886,Q66N60_PELFU/154-232:496.26545
733716114):2.4836610352120747,((((((((((Q9H384_ONCMY/341-533:214.8483993902439,Q7T3U4_9TELE/347-539:217.1516006097561):18.173595369170982,Q762D6
_CONMY/331-524:234.82640463082902):10.848883623633881,ESR2_ICTFO/330-522:251.65111637636613):16.842586741727942,(Q8UW75_ORYLA/314-506:211.836442
13516746,(Q9HT33_ONCMY/338-530:180.31904560810813,(A6FYCE_BOLSO/320-512:168.17313829787233,Q170B4_ORMO/312-504:165.82686170212767):18.6809543
91819793):15.77233195754718,ESR3_MICMU/326-512:201.72276680424528):7.163570648325411):57.524600758272061):44.31278271145291,(((ESR2_CALJA/303-4
94:175.15984513274336,ESR2_MACMU/66-257:172.84015486725664):12.195418552036202,ESR2_PIG/300-490:184.8045414796381):46.5676282051282,Q8MTE6_HUMA
N/303-478:249.6823717948718):74.19112333854709):50.51620536980856,(((A5A588_PARDA/138-309:219.78566576086956,((Q5CCT6_RUTRU/319-511:151.2090336
1344537,Q5NKI4_9TELE/322-514:151.79096638659463):5.021097046413502,Q1HCL4_9TELE/322-514:156.47890295358865):51.71433423913004):25.86272905759132
2,A5A5B7_MISAN/138-309:243.76227094240838):72.74037039620535,(((Q6W5G6_XENLA/184-376:145.42515432098764,Q2SC14_XENTR/342-534:145.5748456790123
6):40.21008522727273,Q9KAL1_SAMPH/341-533:189.289914727273):24.140850360576923,(ALYGA7_EANFA/1-177:212.250998885844747,(((ESR1_PORGO/343-535:150
69037656903765,ESR1_BAY/356-548:153.30962343096235):9.80190677966101,Q6HMS_EUBMA/365-557:161.190893220339):22.249001141552526):23.23414963942
3077):74.0923973695542,(Q5E287_ERFOC/133-325:250.2913445308953,(ESR1_ORENDY/306-457:184.11664746843779,Q6Z8B2_HABBU/307-458:181.883525345222
1):8.064573459715639,(Q800Q2_ZOAVI/311-503:170.98444444444445,Q9DD24_MICSM/356-548:169.01555555555555):24.43542654028436):57.9586550690407):27.4
54477630046593):29.16783272879465):45.96768132194144):62.624510775330236,A7XAMB_ANGAN/1-102:434.87719820904476):25.80073663558082,Q6R765_MONAL/
152-230:480.60160711441915):14.439678808537925):34.12470835639584,Q6U2D7_PIG/1-45:548.2700670342292):43.311292608084459,(((Q8WB79_DROME/310-490:4
    
```

Gambar 18 Output tree dari komputasi *MapReduce*

Program *MapReduce* dengan menggunakan tiga node menghasilkan output yang sama.

Waktu komputasi yang digunakan mapper dalam membangun *phylogenetic tree* berbeda antara Hadoop yang menggunakan dua node dan tiga node, berikut hasil dari tiap percobaan.

Total time spent by all maps in occupied slots (ms)=427104
Total time spent by all reduces in occupied slots (ms)=0
Total time spent by all map tasks (ms)=8898
Total vcore-milliseconds taken by all map tasks=8898
Total megabyte-milliseconds taken by all map tasks=13667328

Gambar 19 Detail waktu yang digunakan komputasi mapeduce dengan dua buah node

Tabel 2 Waktu komputasi menggunakan program *MapReduce*

Jumlah Node	Waktu Komputasi (ms)
2	8898
3	7468

Waktu komputasi yang digunakan Hadoop *MapReduce* dengan dua node dan tiga node berbeda. Hadoop dengan tiga node membutuhkan waktu lebih sedikit dibandingkan dengan Hadoop yang hanya menggunakan dua node.

4.2 Analisis Hasil

Phylogenetic tree dengan IDE Java dan komputasi *MapReduce* menggunakan waktu komputasi yang berbeda, hal ini disebabkan karena perbedaan lingkungan yang dibangun. Waktu komputasi program menggunakan *MapReduce* dengan dua dan tiga buah node hasilnya berbeda. Untuk *MapReduce* yang menggunakan dua node, waktu yang dibutuhkan 8898 ms sedangkan pada tiga node waktu yang dibutuhkan lebih singkat yaitu 7468 ms. Dilihat dari gambar 19, waktu yang digunakan oleh komputasi *MapReduce* diambil dari total waktu yang digunakan *map tasks*.

5. Kesimpulan

5.1 Kesimpulan

Dari hasil pengujian, dapat disimpulkan bahwa:

1. Dalam membangun *phylogenetic tree* dapat menggunakan komputasi *MapReduce*.
2. Dari waktu yang dihasilkan oleh komputasi *MapReduce* dengan menggunakan dua node dan tiga node dapat disimpulkan bahwa jumlah node mempengaruhi waktu komputasi, semakin banyak node maka waktu komputasi semakin cepat.

5.2 Saran

Untuk mengembangkan hasil yang sudah di dapatkan dari penelitian ini, diharapkan pada penelitian selanjutnya dapat:

1. Mendistribusikan *mapper* dan *reducer* secara tepat sehingga waktu komputasi yang digunakan *MapReduce* dapat dipersingkat.
2. Menjalankan secara paralel dan terdistribusi dalam proses iterasi *alignment* yang dibutuhkan untuk melakukan rangkaian *Multiple Sequence Alignment* (MSA).

Daftar Pustaka

- [1] Alberts, Johnson, Lewis, Raff and Roberts, *Molecular Biology of the Cell*, 6th ed., 2014.
- [2] V. O. Polyanovsky, M. A. Roytberg dan V. G. Tumayan, "Comparative Analysis of the Quality of a Global Algorithm and a Local Algorithm for Alignment of Two Sequences," *Algorithms Mol Biol*, 2011.
- [3] R. Chenna, H. Sugawara, T. Koike, R. Lopez, T. J. Gibson, D. G. Higgins dan J. D. Thomp, "Multiple Sequence Alignment with the Clustal Series of Programs," *Nucleic Acids Res*, 2013.
- [4] C. Notredame, D. G. Higgins dan J. Heringa, "T-coffee: a Novel Method for Fast and Accurate Multiple Sequence Alignment," pp. 205-217, 2000.
- [5] J. Alnasir, "The Application of Hadoop in Structural Bioinformatics," *bioRxiv*, 2018.
- [6] C. Ranger, R. Raghuraman, A. Penmetsa, G. Bradski dan C. Kozyrakis, "Evaluating *MapReduce* for Multi-core and Multiprocessor Systems," 2007.
- [7] A. M. Lesk, "Bioinformatics," 26 July 2013. [Online]. Available: <https://www.britannica.com/science/bioinformatics>. [Diakses 1 November 2018].
- [8] Q. Zou, "Multiple Sequence Alignment and Reconstructing Phylogenetic Trees with Hadoop," *2016 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, 2016.
- [9] J. Xiong, *Essential Bioinformatics*, New York: Cambridge University Press, 2006.
- [10] M. Bhandarkar, "MapReduce Programming with Apache Hadoop," *2010 IEEE International Symposium on Parallel & Distributed Processing (IPDPS)*, 2010.
- [11] A. Lafita, S. Bliven, A. Prlic, D. Guzenko, P. W. Rose, A. Bradley, P. Pavan, D. Myers-Turnbull, Y. Valasatava, M. Heuer, M. Larson, S. K. Burley dan J. M. Duarte, "BioJava 5: A Community Driven Open-source Bioinformatics Library," *PLoS Comput Biol* 15(2): e1006791, 2019.