

Abstract

Similarity searching tools are generally divided to two distinct groups: Alignment-based and alignment-free methods. There is no standard comparison between alignment-based and alignment-free methods` performance. In this project, I compared CLUSTALW, T-COFFEE and MUSCLE algorithms with four alignment-free methods, Kr, D2, D2z and MplusD to find out which of these methods are more efficient with respect to time and accuracy. In both alignment-based and alignment-free, results are represented by phylogenetic trees. Here, all phylogenetic trees are compared against the reference tree to determine how accurate the methods are.

I found that alignment-based methods were the most accurate comparing to alignment-free methods. On the other hand, Kr and D2z gave approximately correct results with substantial gain in speed.

I conclude that for inferring phylogenetic trees containing taxa that have diverged for a long time, alignment-based methods should be used. However if many closely related species are studied, then an alignment-free method should worth well.

Abstrak

Perkakas perbandingan keserupaan boleh dibahagikan secara am kepada dua kumpulan: kaedah berasaskan penjujukan dan kaedah bebas penjujukan. Walaupun kaedah berasaskan penjujukan lebih diutamakan kerana ia berkisar pada konsep homologi, ia biasanya mengambil masa yang lama apabila kita ingin menjujukan turutan-turutan yang sangat panjang daripada banyak taxa. Sementara itu, kaedah bebas penjujukan yang berasaskan frekuensi perkataan dalam suatu jujukan lebih cepat secara relatif dan mungkin memberikan topologi yang lebih kurang tepat. Dalam projek ini, saya telah membandingkan tiga kaedah berasaskan penjujukan: CLUSTALW, T-COFFEE dan MUSCLE terhadap empat kaedah bebas penjujukan: Kr, D2, D2z dan MplusD untuk mengetahui yang manakah kaedah ini lebih cekap untuk membina pepohon filogenetik dalam aspek kelajuan dan ketepatan, dengan menggunakan algoritma *Neighbour Joining*. Dengan menggunakan suatu pohon rujukan yang mengandungi *clade* primat dan burung yang dibina daripada jujukan genom mitokondria dan dijustifikasikan secara biologi, saya mendapati bahawa kaedah berasaskan penjujukan adalah paling tepat berbanding kaedah bebas penjujukan – mereka semua memberikan topologi yang sama seperti pohon rujukan. Sementara itu, Kr dan D2z memberikan keputusan yang lebih kurang tepat dengan kelebihan kelajuan. Saya merumuskan bahawa untuk menganggarkan pepohon filogenetik yang mengandungi taxa yang sudah berpisah untuk masa yang lama, kaedah berasaskan penjujukan harus dipakai. Walaubagaimanapun, jika banyak spesis yang berkait rapat dikaji, maka barangkali kaedah bebas penjujukan agak berguna.

Acknowledgment

I am really grateful to my supervisor Dr. Tsung Fei Khang for giving me the opportunity to involve in this project and devoting his valuable time and unlimited support for directing me all the way throughout this project. I want to take this opportunity to express my gratitude to him for not only guiding me in this project but also helping me to open my mind into new area of feature technology.

I would like to thank my family members, especially my husband, Masoud, for supporting and encouraging me to pursue my project.

Last but not least, my thanks and appreciations go to my friend Hoda for helping me every time and everywhere, when I was facing difficulty.

Table of Contents

Abstract	i
Abstrak	ii
Acknowledgment	iii
Table of Contents	iv
List of Figures	vi
List of Tables.....	vi
1.0 Chapter 1: Introduction	1
1.1 Sequence Comparison	1
1.2 Structure of Mitochondrial Genome	5
1.3 Research Objectives	6
2.0 Chapter 2 Literature Review	7
2.1 Alignment-based methods	7
2.1.1 CLUSTALW	8
2.1.2 MUSCLE	10
2.1.3 T-COFFEE	13
2.1.4 Limitations of Alignment-Based methods	17
2.2 Alignment-free methods.....	18
2.2.1 Kr	19
2.2.2 D2.....	20
2.2.3 D2z	21

2.2.4 MplusD.....	22
2.2.5 Limitations of Alignment-Free methods.....	23
2.3 Sequences GC content	24
3.0 Chapter 3 Methods and Materials	25
3.1 Collection of Data	25
3.2 Software and Tools	25
3.3 Genome Comparison.....	27
4.0 Results	30
4.1 GC Content Sliding Window Plot.....	30
4.2 Phylogenetic trees	32
5.0 Discussion	36
6.0 Conclusions	38
References	39

List of Figures

Figure 1.1. Data Flow	4
Figure 1.2. The Human mitochondrial genome.	5
Figure 2.3. CLUSTALW workflow	9
Figure 2.4. The flow of MUSCLE algorithm.....	11
Figure 2.5. Basic procedure of T-Coffee	14
Figure 2.6. Mitochondrial genome comparison among primates..	27
Figure 2.7. Mitochondrial genome comparison among birds.	28
Figure 2.8. Mitochondrial genome comparison between human and chicken.....	28
Figure 2.9. Mitochondrial genome comparison between human and brachiopoda	29
Figure 2.10. Mitochondrial genome comparison between chicken and brachiopoda.....	29
Figure 4.11. GC content sliding window pot.	30

List of Tables

Table 4.1. Phylogeny trees resulted from alignment-based methods.....	32
Table 4.2. Phylogeny trees resulted from alignment-free methods	34
Table 4.3. Comparison of phylogeny trees	35