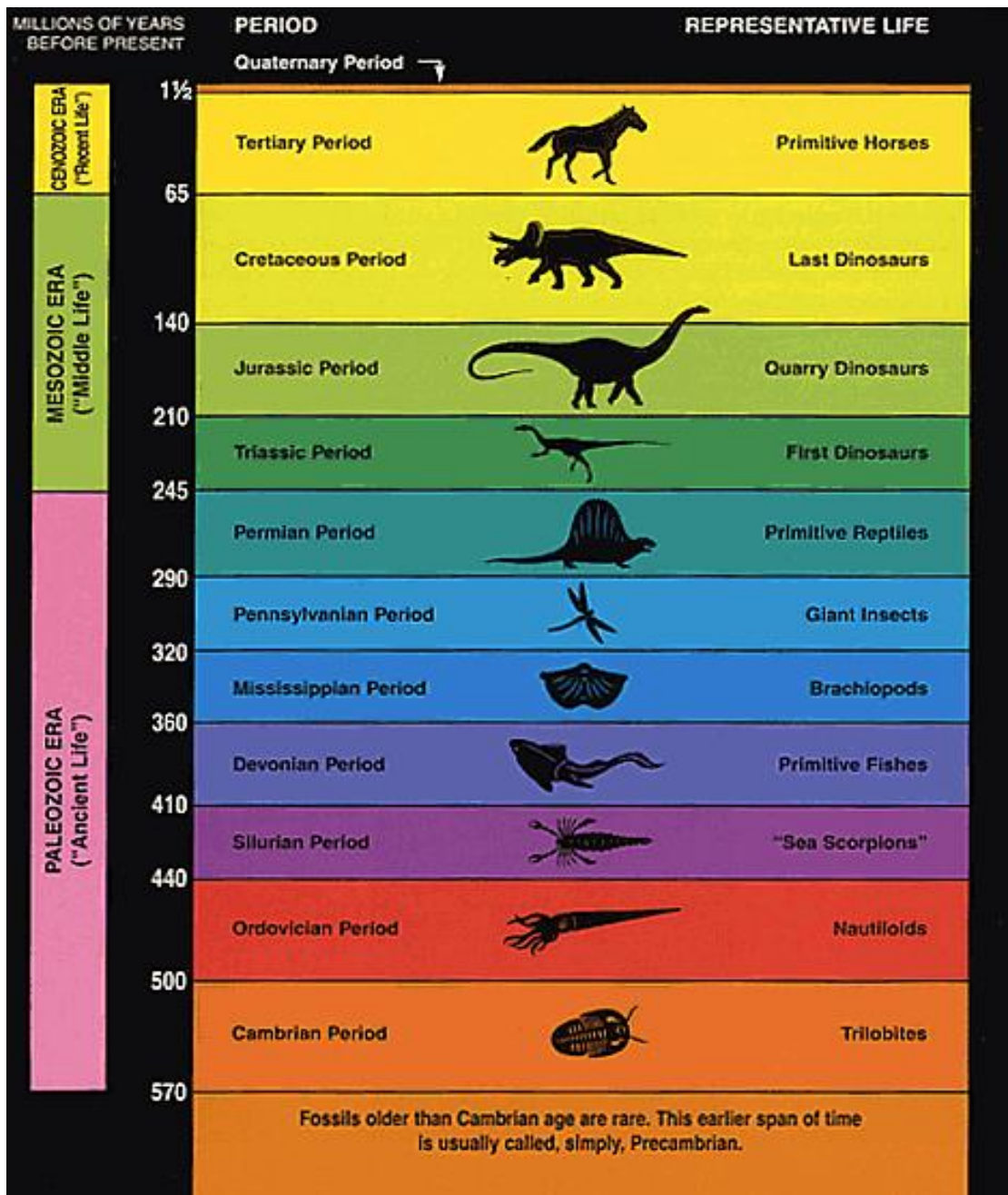


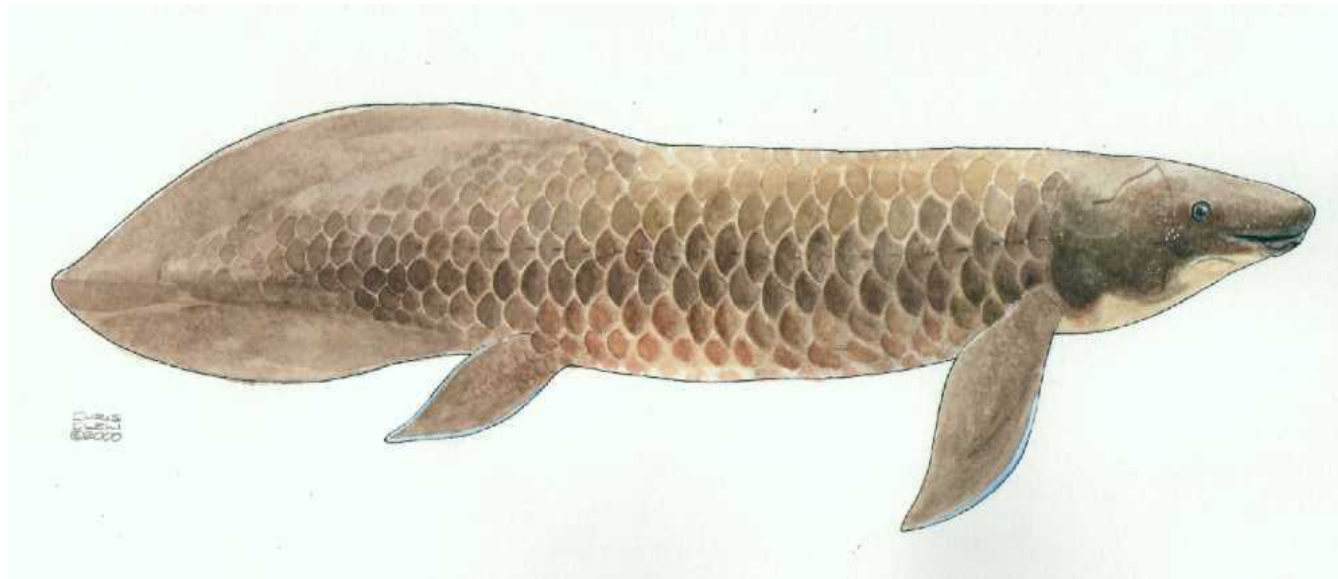
Statistical Learning Based on Distributions of Oligonucleotides in DNA Sequences

Probal Chaudhuri

Indian Statistical Institute, Calcutta



Australian Lungfish



Coelacanth, the living fossil



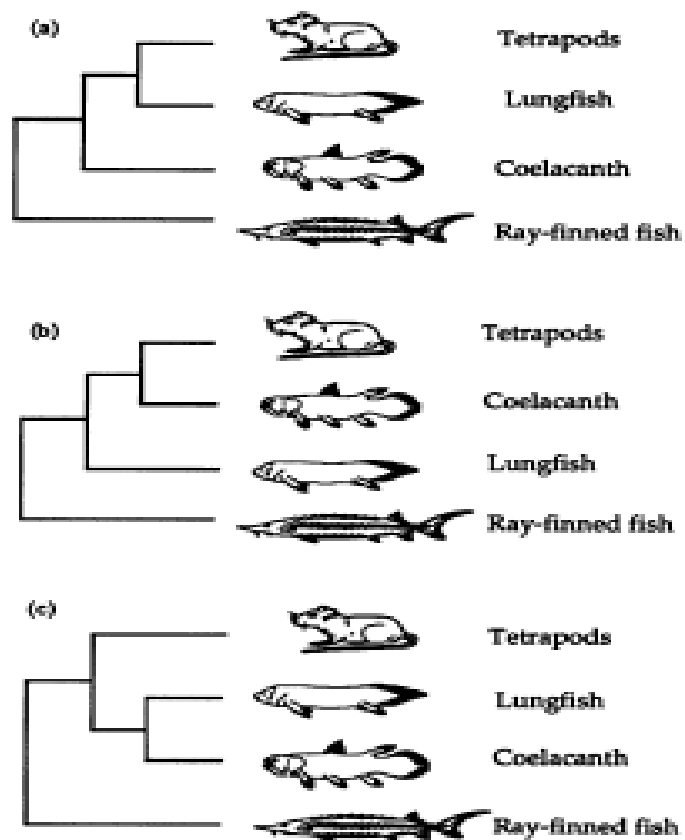


FIG. 1. Alternative hypotheses of sister group relationships between sarcopterygii and tetrapods. (A) Lungfish as the sister group of tetrapods. (B) Coelacanth as the closest living relative of tetrapods. (C) Coelacanth and lungfish equally closely related as sister groups of tetrapods.

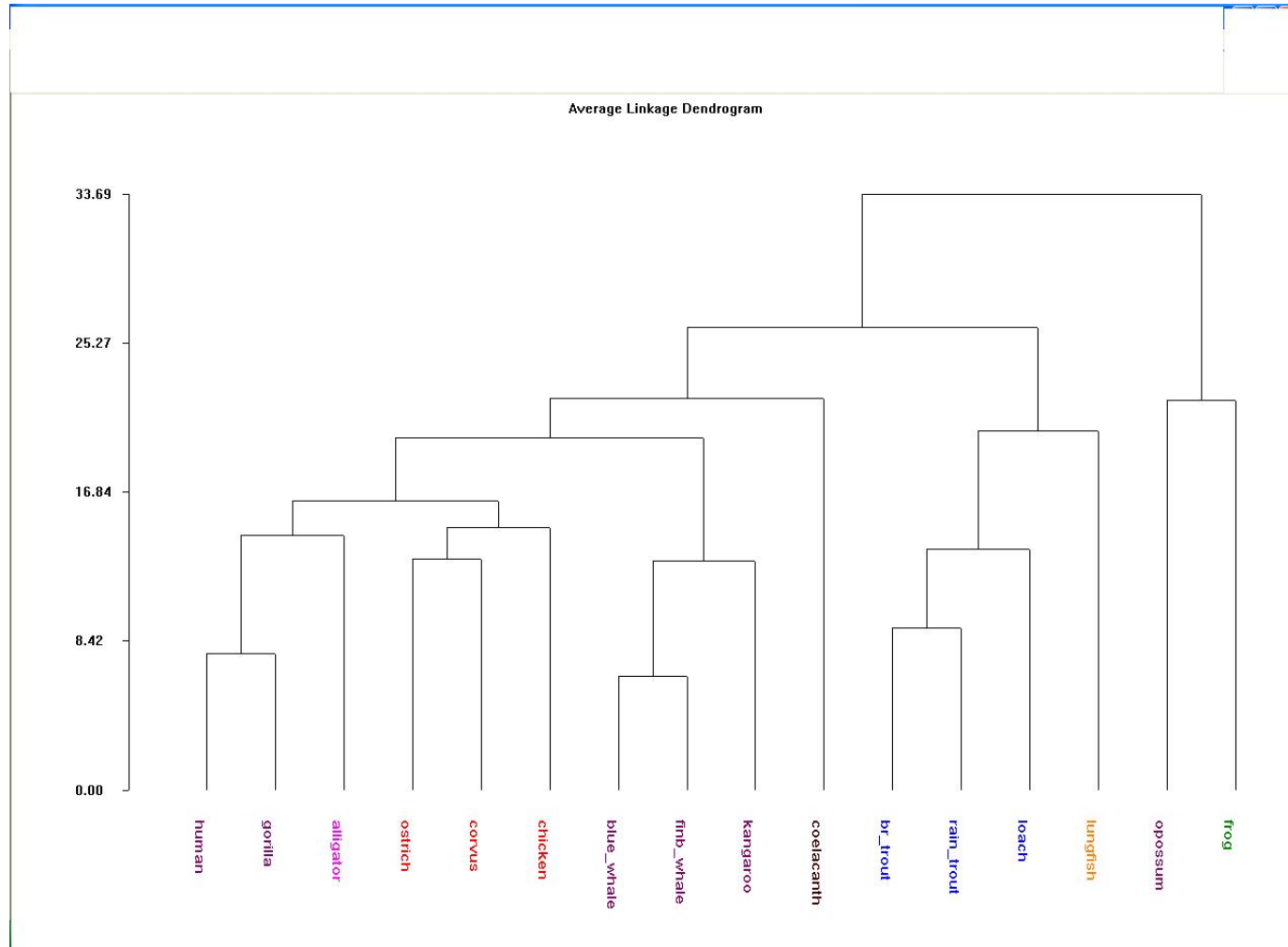
Analysis of Mitochondrial Genomes

- DNA sequences of complete mitochondria of several vertebrates.
- Lengths of sequences varying from 16K to 18K nucleotide bases.
- The sample includes a few mammals, a few birds and a few standard types of fish. It also contains a reptile and an amphibian animal.

Statistical Analysis of DNA Sequences

- Oligo-nucleotide frequencies are computed for each mitochondrial sequence.
- These frequencies can be used as numerical measurements or variables. They form multivariate observations associated with the given DNA sequences.
- One can carry out cluster analysis with this multivariate data.

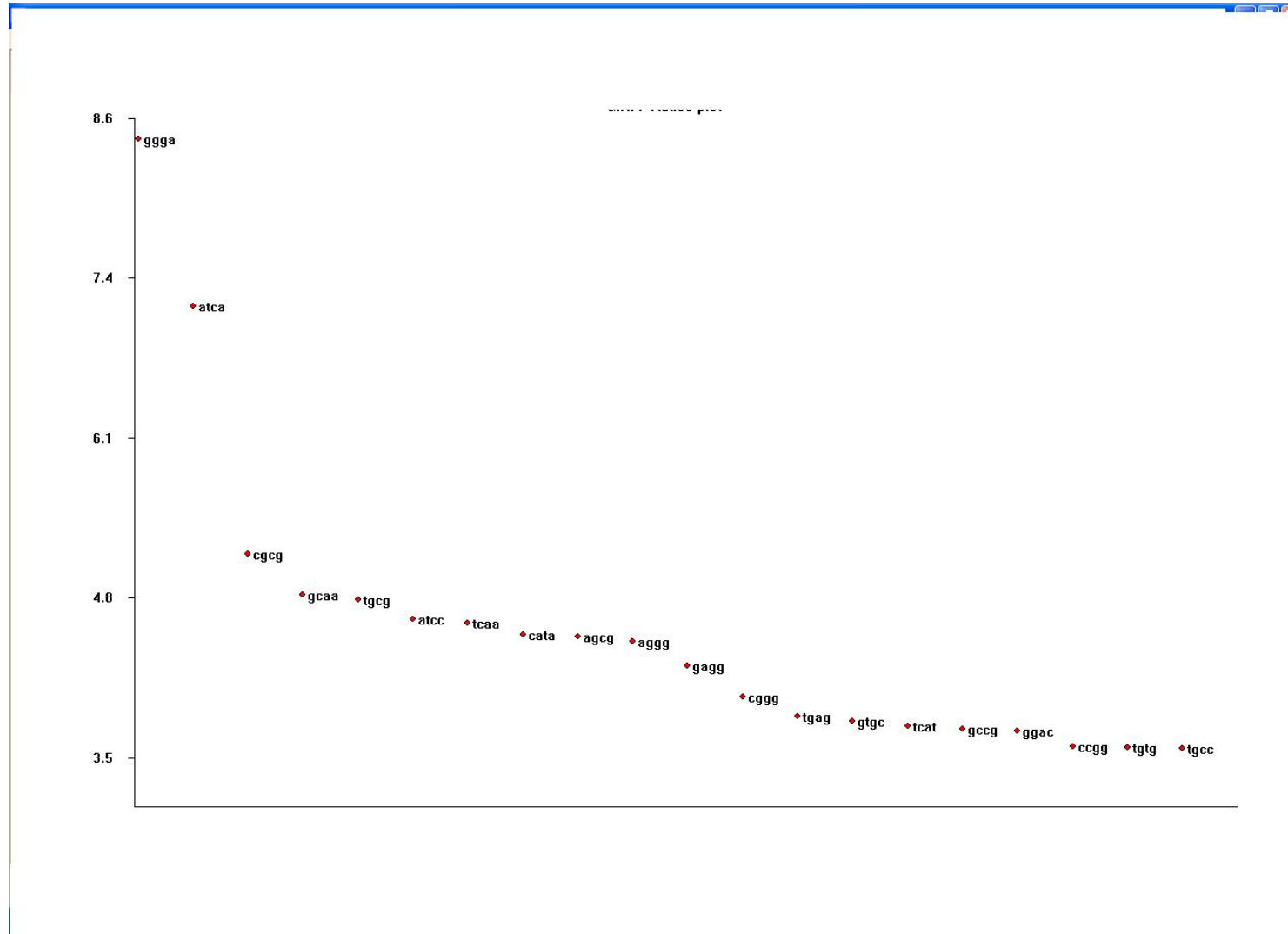
Result of cluster Analysis



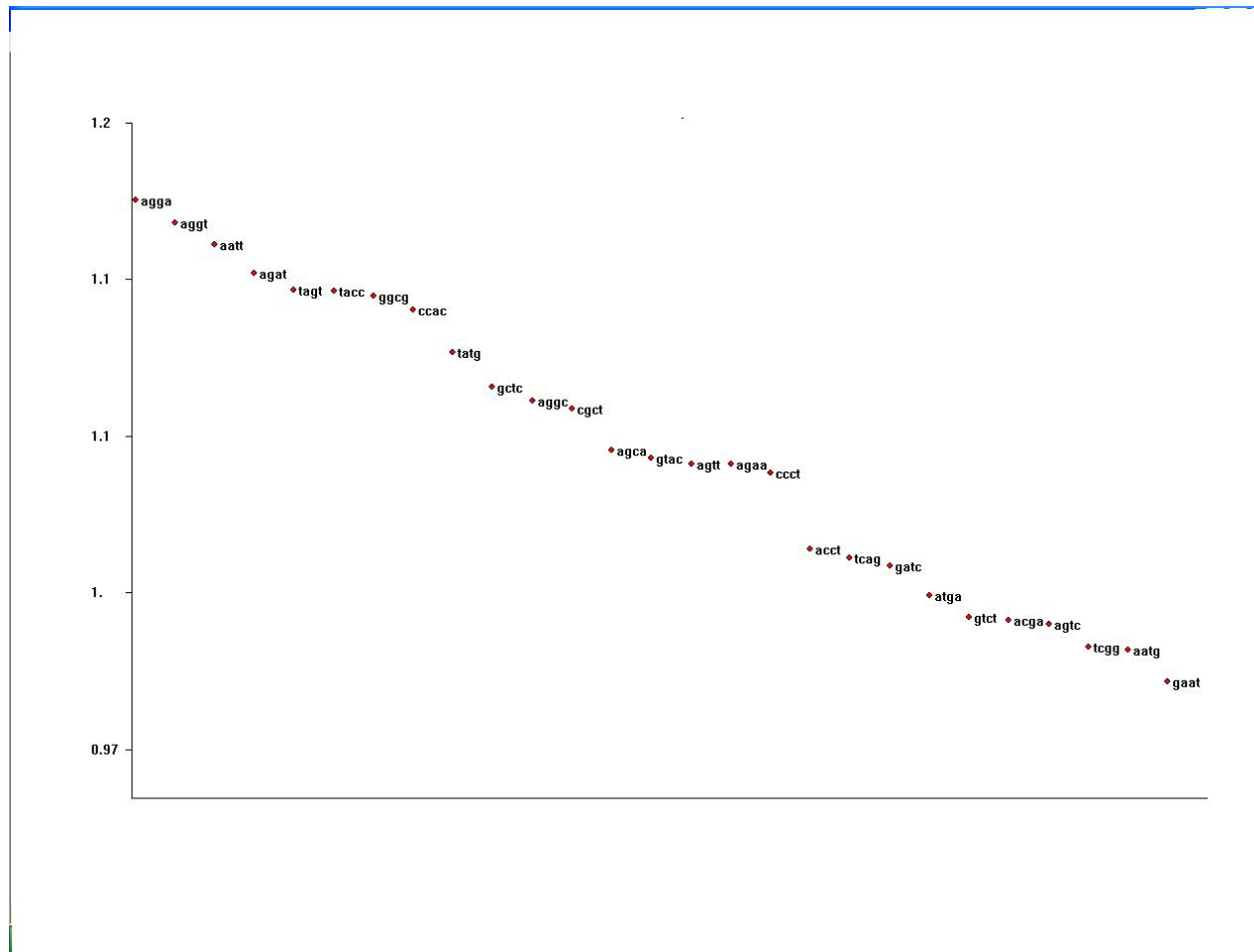
Identification of Informative Oligonucleotides

- The sample contains mammals, birds and standard fish, which form three distinct groups.
- Compute within group and between group variations of each oligonucleotide frequency.
- Compute the ratio of those two variations for each oligonucleotide.
- Plot the ratios for different oligonucleotides.

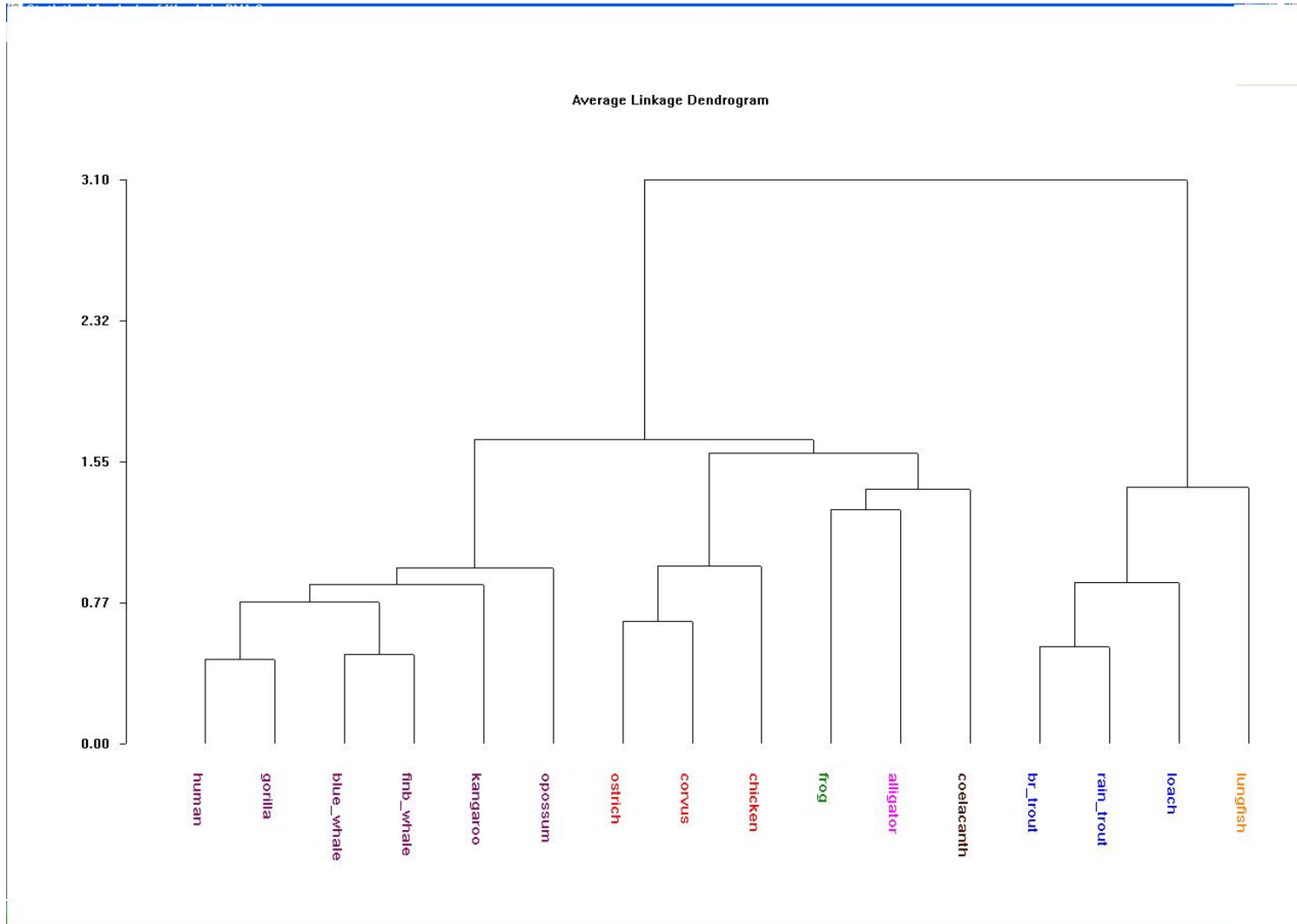
Oligonucleotides with Large Ratios



Oligonucleotides with Small Ratios



Refined Cluster Analysis



Identification of Genomic Islands

- Genomic Islands (GIs) are those parts of a bacterial genome that have originated through horizontal gene transfer.
- In pathogenic bacteria some of the GIs are often associated with pathogenic activities.
- GIs may also contain genes that are responsible for symbiotic activities and may help the bacteria in the process of adaptation.

Since GIs are acquired by the bacterial genome by horizontal transfer, one may expect that it will have oligonucleotide composition different from that in other parts of the genome.

Fig.1

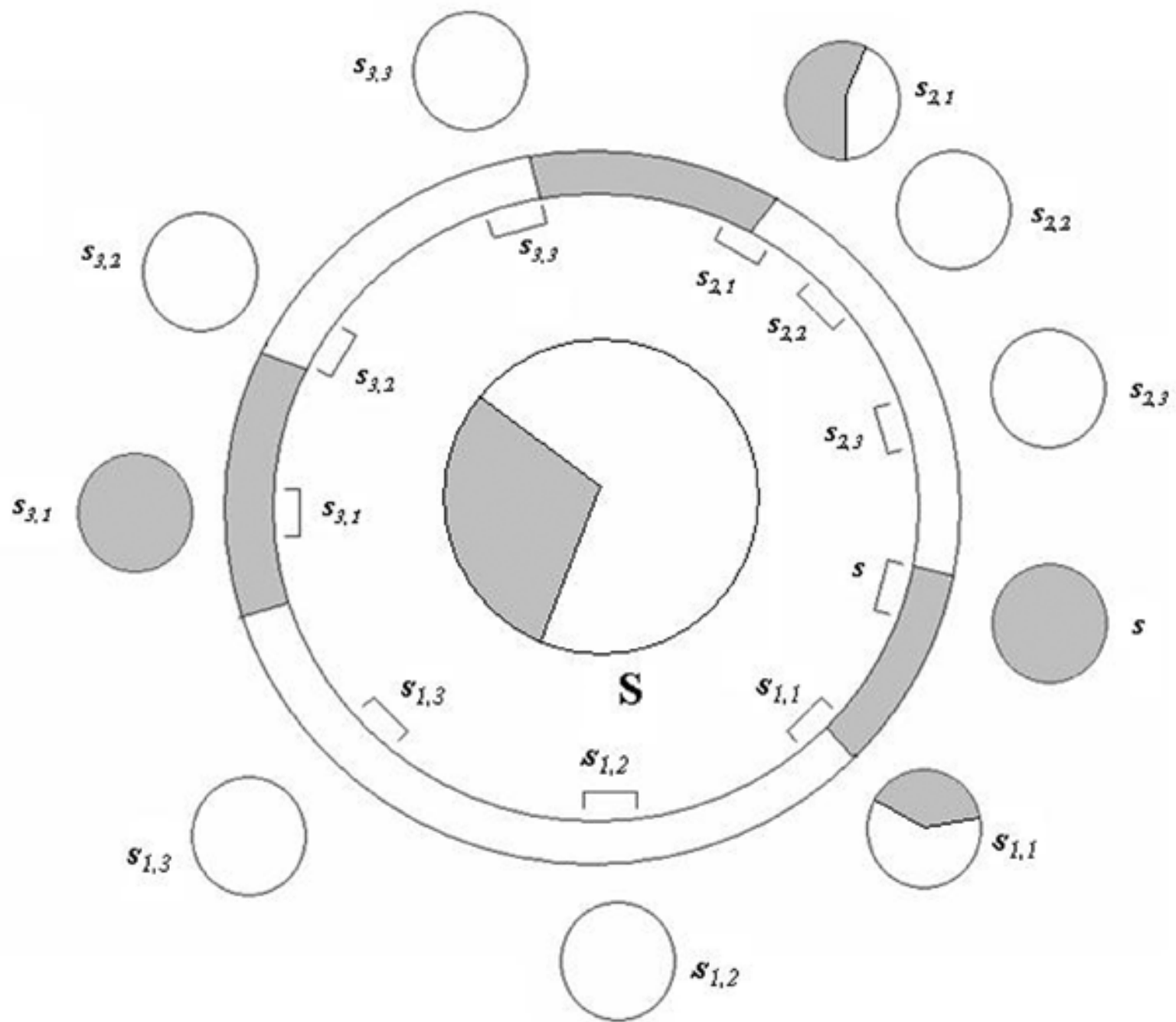


Fig.3A

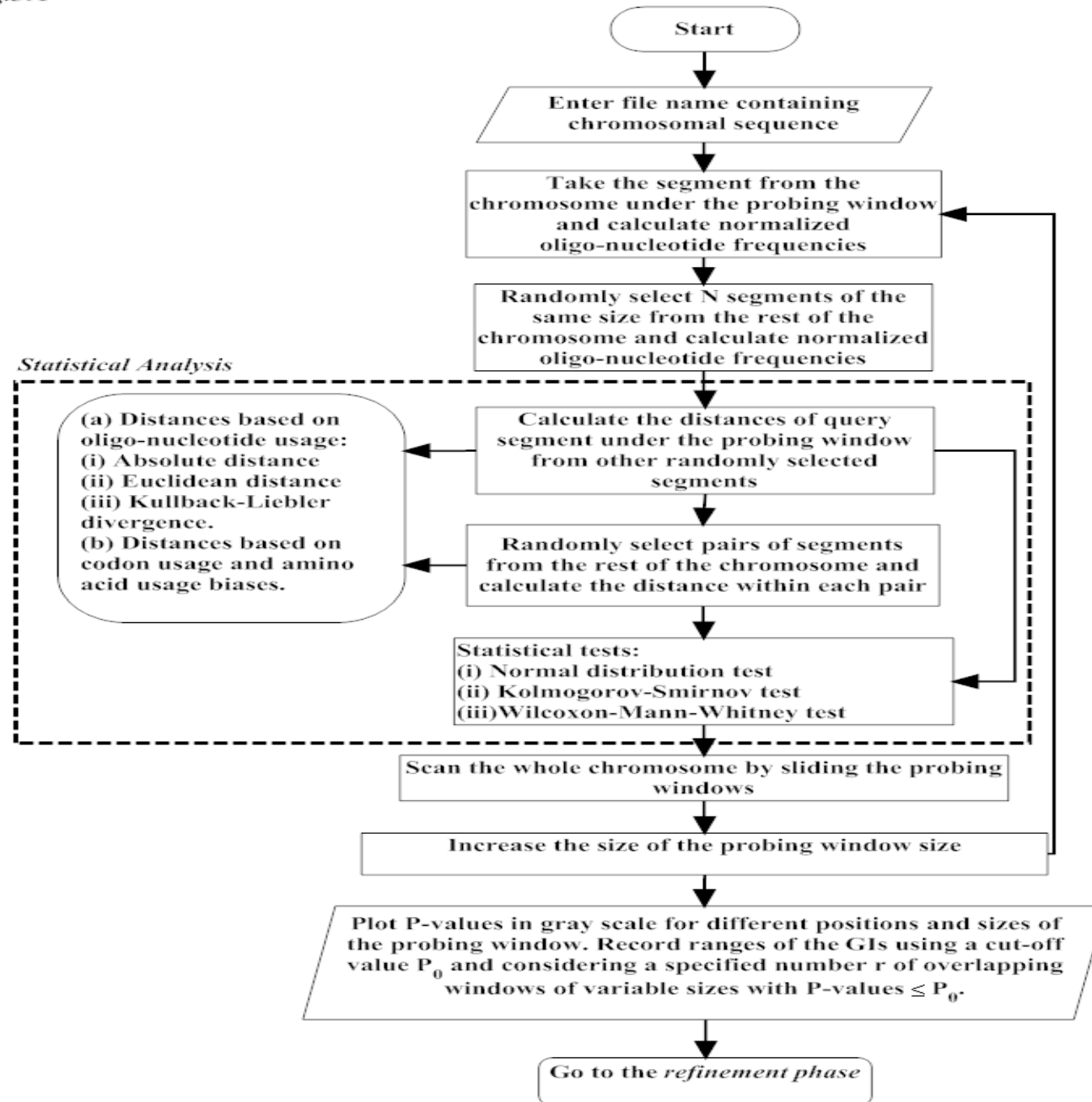
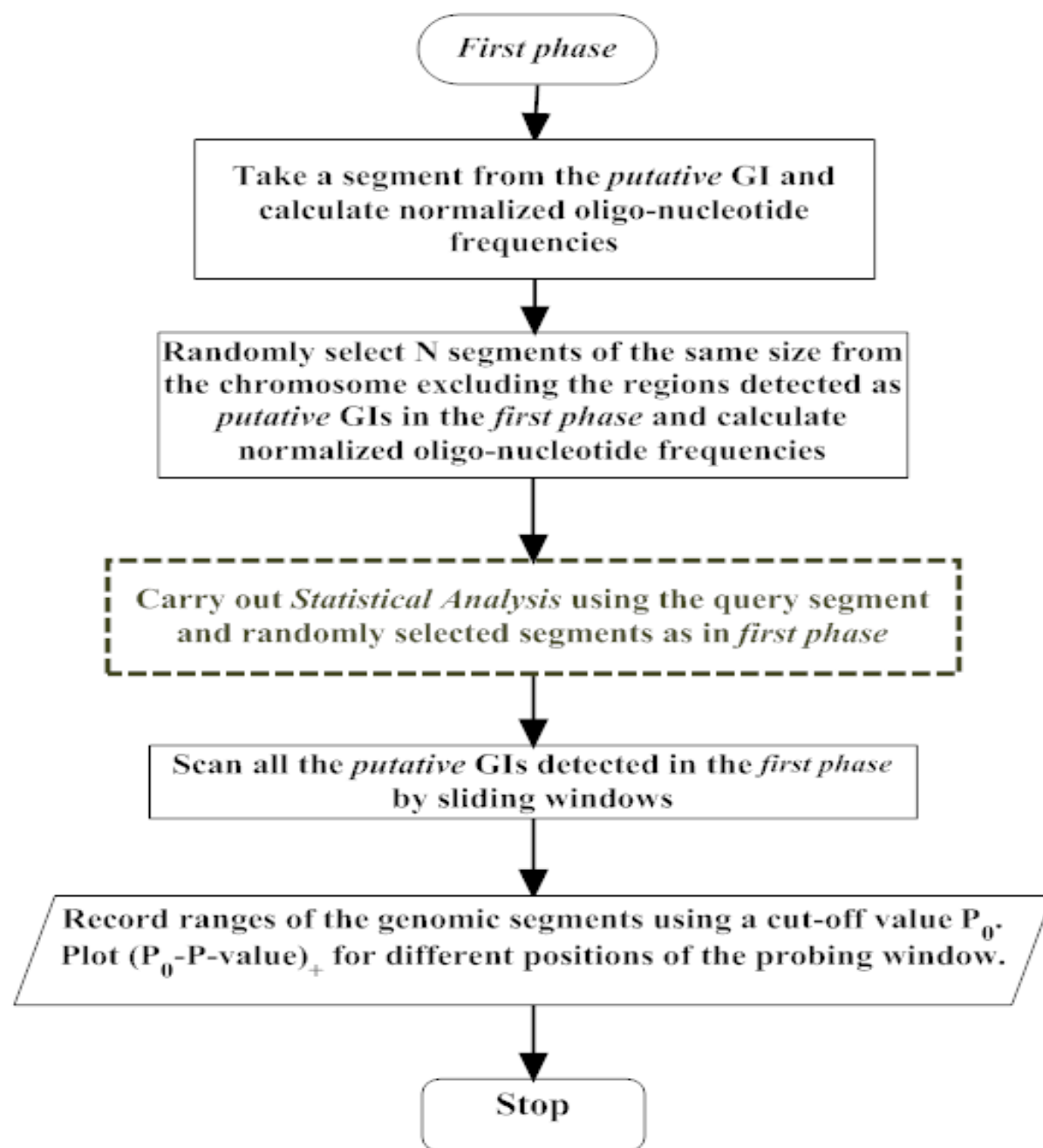


Fig.3B



Salmonella typhi CT18

- A pathogenic bacteria causing typhoid.
- The genome is known to contain several pathogenic islands.
- Many of these islands are reported to contain genes from bacteriophages.

Fig.4A(i)

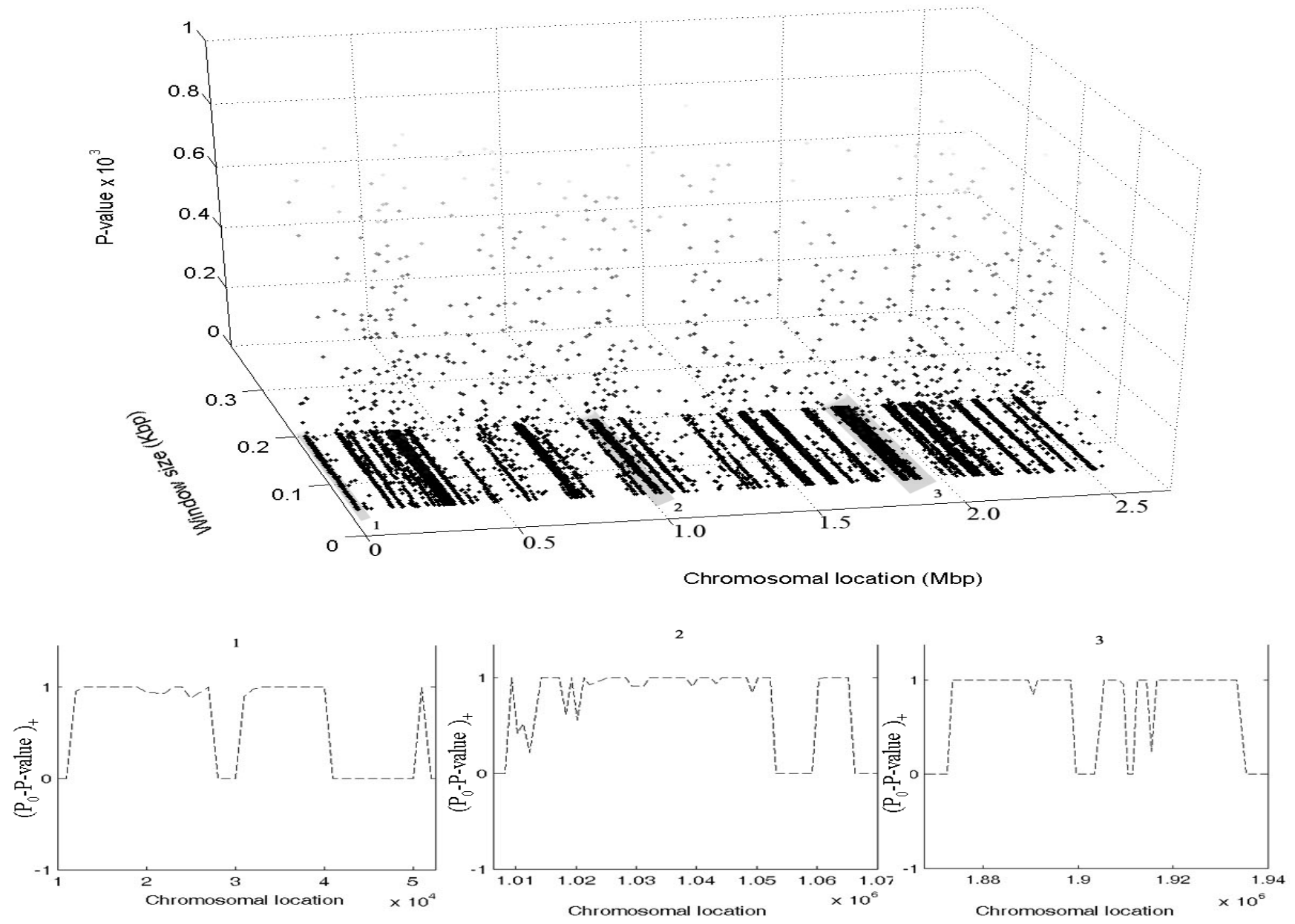
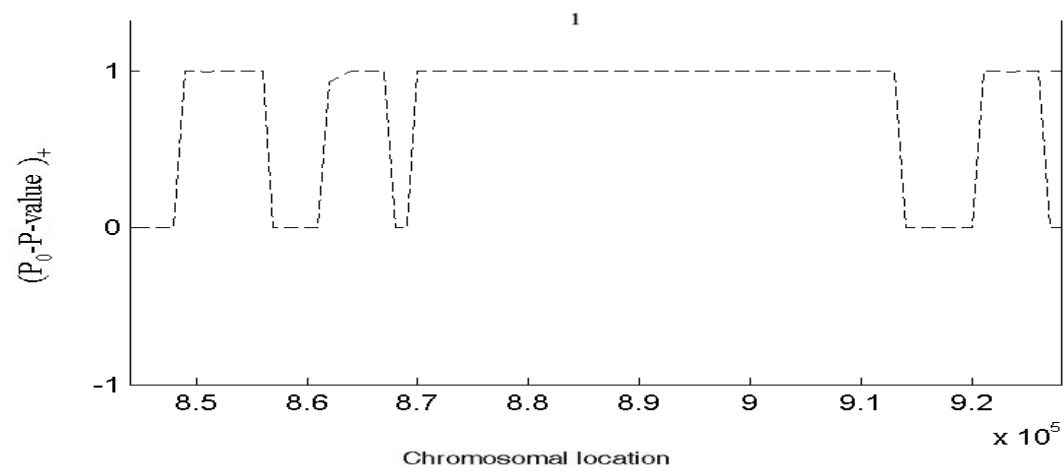
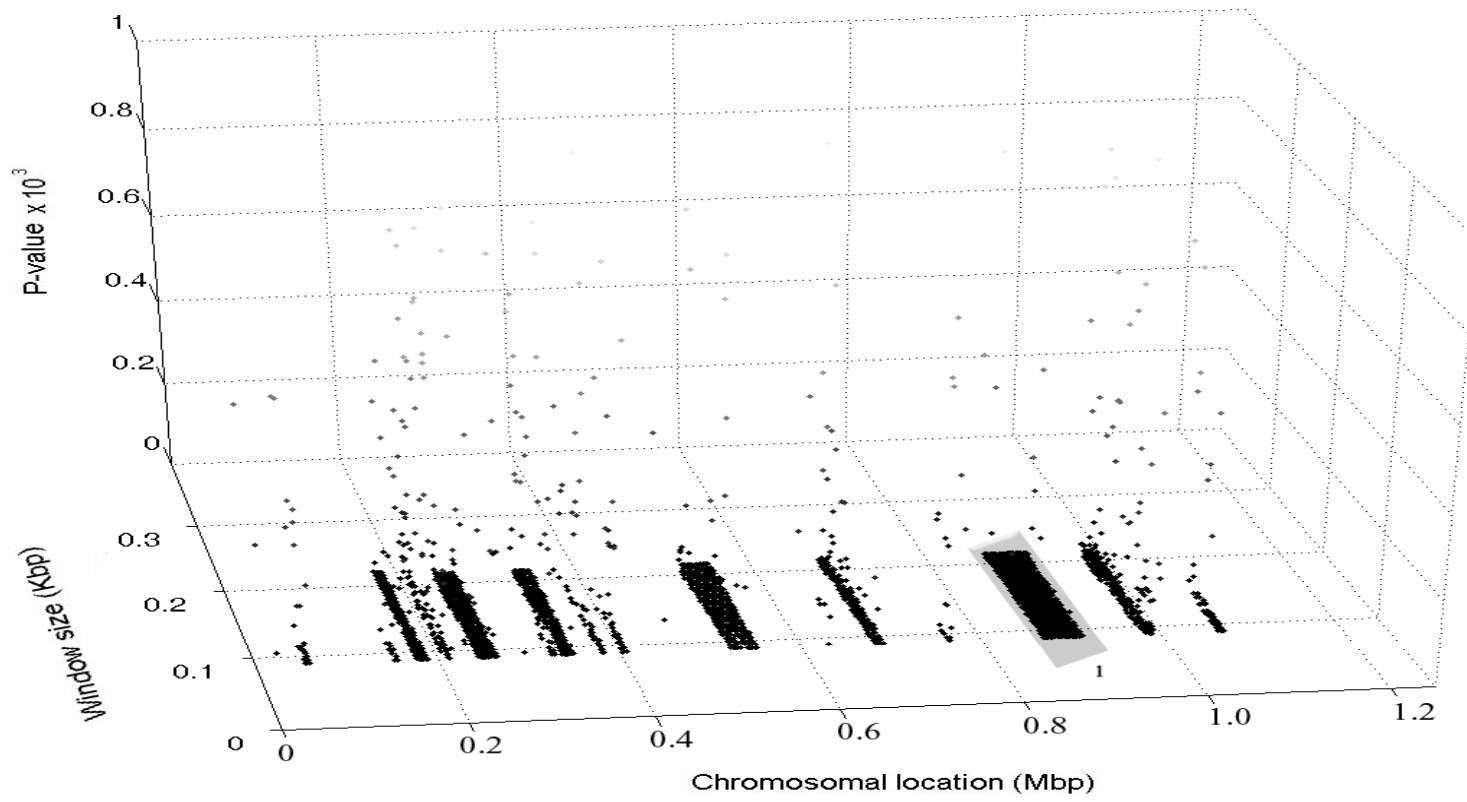


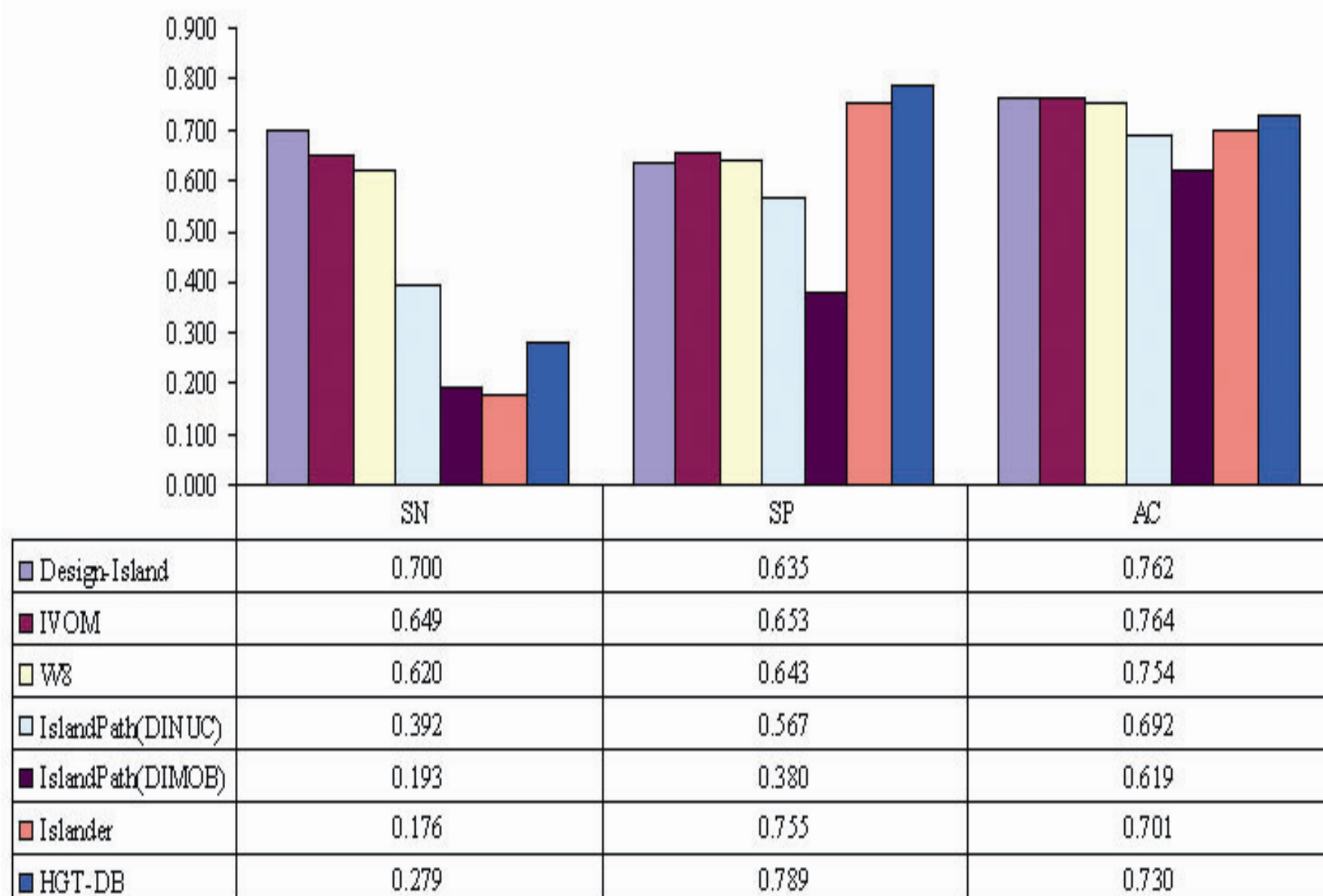
Fig.4B(i)



Salmonella typhi CT18 (contd.)

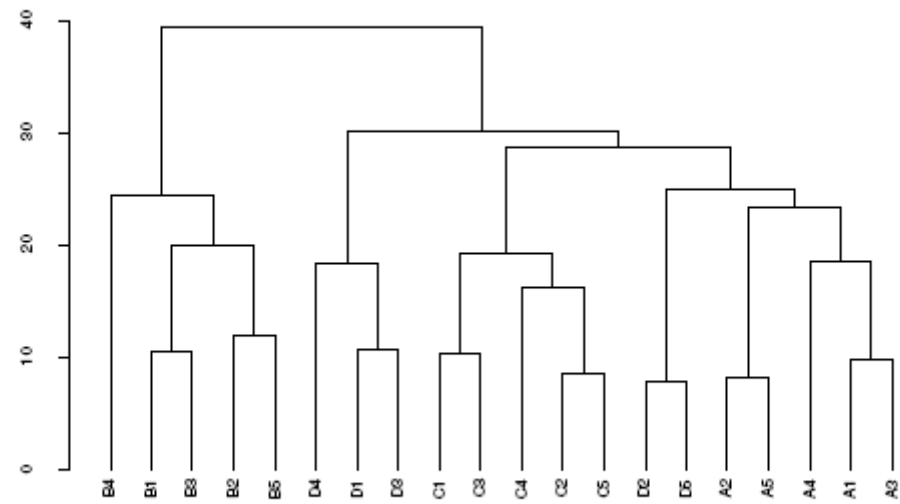
- *S. typhimurium* LT2 is considered a sister lineage to *S. typhi* CT18 and *E. coli* K12 is considered an outgroup (Vernikos and Parkhill, 2006, *Bioinformatics*).
- Genes present in all three genomes form a set of Core Genes.
- A data set consisting of 1560 manually curated **putative horizontally transferred genes**.

Fig.5

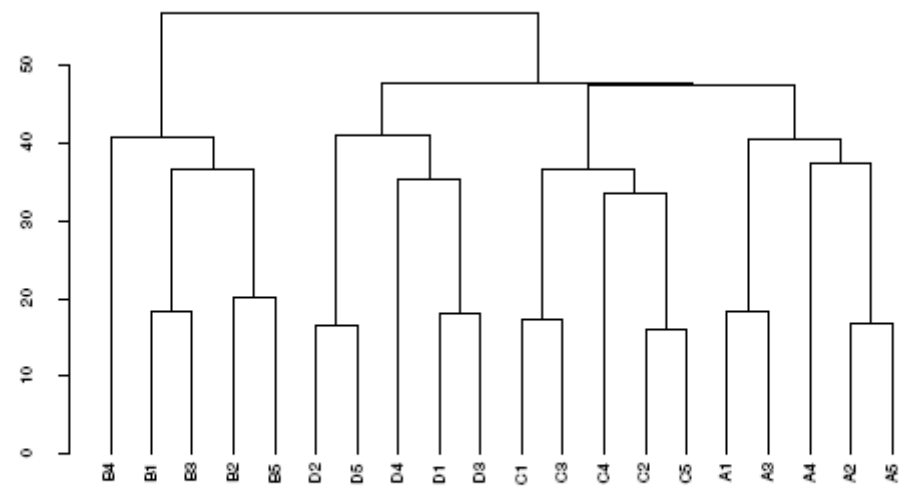
Salmonella typhi CT18

Data on bacteriophage genomes

- Four different single stranded lytic phages :
ΦX174 (host : *E. coli*), *G4* (host : *E. coli*),
F1 (host : *E. coli*) and *PF3* (host : *P. aeruginosa*)
- The sizes of these phage genomes lie between 5.3K to 6.5K nucleotide bases.
- We consider five segments of each of these phage genomes, and carry out cluster analysis based on oligonucleotide frequencies.

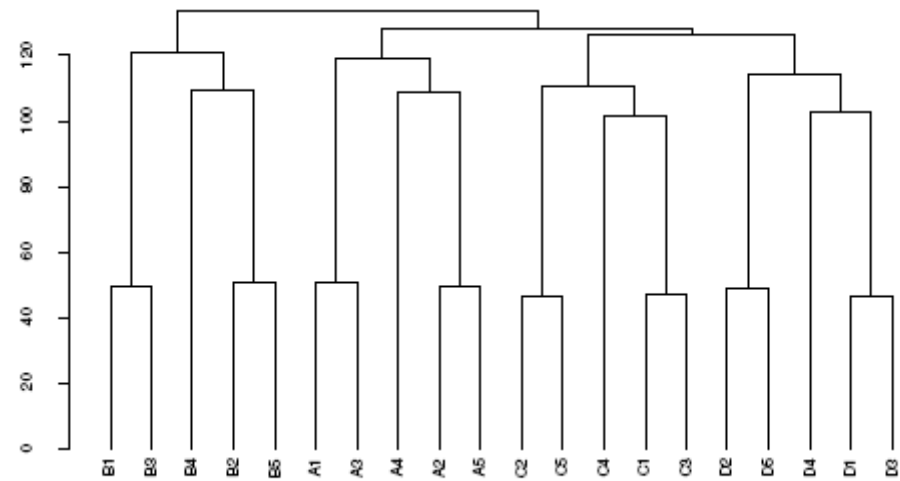


3-word frequencies

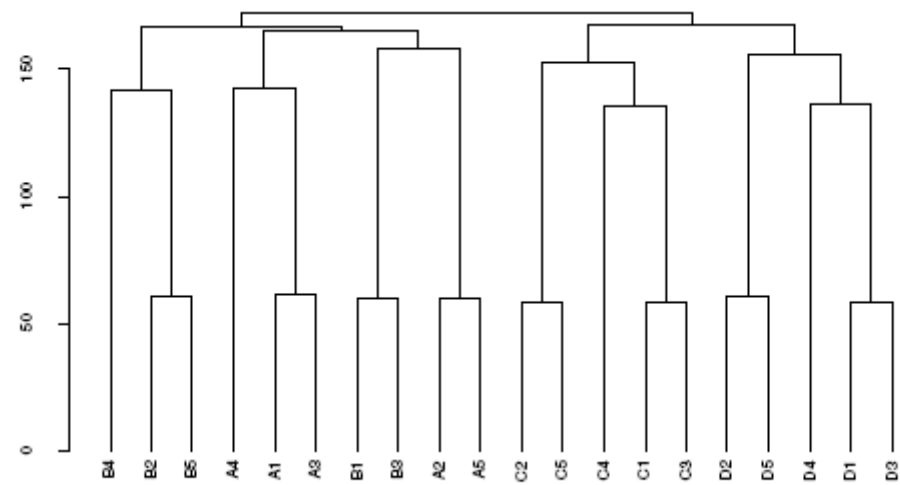


4-word frequencies

Fig. 2.5. Dendrograms for average linkage cluster analysis of fragments of phage genomes.



Dendrogram for 6-word frequencies



Dendrogram for 7-word frequencies

Fig. 2.6. Average linkage cluster analysis of fragments of bacteriophage genomes.

References

- P. Chaudhuri and S. Das (*Current Science*, 2001)
- S. Basu, D. P. Burma and P. Chaudhuri (*Journal of Mathematical Biology*, 2003)
- G. S. Vernikos and J. Parkhill (*Bioinformatics*, 2006)
- S. Nag, R. Chatterjee, K. Chaudhuri and P. Chaudhuri (*Sadhana*, 2006)
- R. Chatterjee, K. Chaudhuri and P. Chaudhuri (*BMC Genomics*, 2008)