

# Whole Genome Phylogeny Based on Clustered Signature String Composition

Xiaomeng Wu, Guohui Lin\*  
Department of Computing Science  
University of Alberta  
Edmonton, Alberta T6G 2E8, Canada  
xiaomeng,ghlin@cs.ualberta.ca

Xiu-Feng Wan, Dong Xu  
Department of Computer Science  
University of Missouri – Columbia  
Columbia, Missouri 65211, USA  
wanx,xudong@missouri.edu

## Abstract

*Peptide compositions constructed out of whole sets of protein sequences can be used as species signatures for phylogenetic analysis. To account for point mutations, an amino acid substitution model is integrated into the complete composition vectors through a novel peptide clustering algorithm. Such a refined signature is expected to highlight deeper evolutionary relationships among the species and employed into the whole genome phylogenetic analysis to define a new evolutionary distance measure. Computational experiments have been set up to validate the effectiveness of this new measure and a vertebrate evolutionary tree using a dataset of 832 proteins for 64 vertebrates is reported.*

## 1. Introduction

The availability of an increasing number of completely sequenced genomes has opened up new avenues for understanding the evolution. In contrast to the traditional approaches where the molecular data is usually cautiously selected, the whole genomes afford unprecedented opportunities and perspectives for detecting evolutionary relationships at a micro point of view. However, this vast amount of sequence data challenge the phylogenetic analyses for evolutionary information representation to digest molecular sequences of millions of bytes.

The carefully selected data in the traditional approaches is relatively easy to be analyzed by adopting some substitution models that describe the prior knowledge about the evolution model. This seems to be advantageous over whole genomes in which gene transfer, unrecognized paralogy, and highly variable evolution rates exist. However, though the phylogenetic analysis by traditional approaches could provide accurate results on the selected molecular data, it

doesn't tell well the species evolution since different sets of selected data normally result in conflicting phylogenetic analyses. On the other hand, whole genomes are believed to contain the complete evolutionary information and the phylogenetic analyses based on whole genomes are expected to equate the evolution of the organisms. Therefore, whole genome phylogeny becomes one of the major problems in comparative genomics [1]. The most profound difficulty in building phylogenies using whole genomes is to effectively and efficiently represent the evolutionary information hidden in whole genomes.

Traditional character-based phylogeny construction methods, including Maximum Parsimony (MP) and Maximum Likelihood (ML), build trees that optimize the distribution of the molecular data for each character, where substitution models are taken to align the multiple entries. Whole genomes for different organisms might contain different sets of genes in different sequential orders on the chromosomes. Therefore, multiple alignment can no longer be applied, not to mention its high computational complexity. During the past a few years, a considerable amount of efforts have been devoted to whole genome phylogeny study. All these efforts successfully avoid the high complexity stage of multiple alignment, and try to use the evolutionary information hidden in whole genomes as much as possible without using a substitution model. The main difference among these efforts is how they treat whole genomes to define the relative evolutionary distance between two whole genomes. Once the pairwise evolutionary distance matrix for the set of taxa is computed, they subsequently call distance-based phylogeny construction methods such as Neighbor-Joining [4] to build the tree.

Among several approaches, one category of whole genome phylogeny methods use the frequencies of segments of amino acids (or nucleotides if DNA sequences) as the species signatures. For example, frequencies of segments of all possible lengths are included in the *complete information set* [2]; linear combinations of frequencies of tri/tetra-peptides using a singular value decomposition are

\*To whom correspondence should be addressed.

