

Identification of Ortholog Groups in KEGG/SSDB by Considering Domain Structures

Masumi Itoh

itoh@kuicr.kyoto-u.ac.jp

Akihiro Nakaya

nakaya@kuicr.kyoto-u.ac.jp

Minoru Kanehisa

kanehisa@kuicr.kyoto-u.ac.jp

Bioinformatics Center, Institute for Chemical Research, Kyoto University, Gokasho, Uji, Kyoto 611-0011, Japan

Keywords: ortholog, clustering, domain extraction

1 Introduction

Huge amount of genome information is stored in databases with the advent of recent genome projects. Although we can effectively predict protein sequences from these genomes, functions of most proteins are not experimentally determined. Therefore computational methods are most important for the function prediction, based on comparison and clustering of protein sequences. However, complications arise from the fact that the unit of conservation is not entire protein molecules but domains which are parts of the protein molecule. Hence a method to classify proteins according to their domain structures must be developed for use in functional predictions. Here, we propose a method for extracting domain information from a cluster of similar proteins obtained by all to all pairwise sequence comparisons of completely sequenced genomes.

2 Material and Method

The KEGG/SSDB database contains Smith-Waterman similarity scores of about 100,000,000 pairs from 350,000 proteins in 100 genomes of KEGG/GENES [2]. Our method performs domain extraction and fine protein clustering for a given group of similar proteins by the following procedures.

1. Construction of similarity profiles for each residues:

Extract one protein (target protein) from the group and compare against all other proteins. Construct a bit vector for each residue of the target protein, in which the bit is one if the residue is aligned to the corresponding protein or zero if not aligned (Fig. 1).

2. Self-comparison of similarity profiles

Calculate similarity scores, which is defined by the equation in bottom of Figure 1, for all amino acid pairs within the target protein (Fig. 2).

3. Extraction of domains

Detect the position where the similarity score between the current position and the next position becomes higher than the similarity score between the current position and the preceding position along the amino acid sequence. This is considered as a boundary of domain candidates.

