

BETAWARE: a machine-learning tool to detect and predict transmembrane beta-barrel proteins in prokaryotes

Castrense Savojardo^{1,2}, Piero Fariselli^{1,2,*} and Rita Casadio¹¹Biocomputing Group, CIRI-Health Science and Technology/Department of Biology, University of Bologna, 40126 Bologna and ²Department of Computer Science and Engineering, University of Bologna, 40127 Bologna, Italy

Associate Editor: Anna Tramontano

ABSTRACT

Summary: The annotation of membrane proteins in proteomes is an important problem of Computational Biology, especially after the development of high-throughput techniques that allow fast and efficient genome sequencing. Among membrane proteins, transmembrane β -barrels (TMBBs) are poorly represented in the database of protein structures (PDB) and difficult to identify with experimental approaches. They are, however, extremely important, playing key roles in several cell functions and bacterial pathogenicity. TMBBs are included in the lipid bilayer with a β -barrel structure and are presently found in the outer membranes of Gram-negative bacteria, mitochondria and chloroplasts. Recently, we developed two top-performing methods based on machine-learning approaches to tackle both the detection of TMBBs in sets of proteins and the prediction of their topology. Here, we present our BETAWARE program that includes both approaches and can run as a standalone program on a linux-based computer to easily address in-home massive protein annotation or filtering.

Availability and implementation: <http://www.biocomp.unibo.it/~savojard/betawarecl>

Contact: piero.fariselli@unibo.it

Supplementary information: Supplementary data are available at *Bioinformatics* online.

Received on October 26, 2012; revised on December 13, 2012; accepted on December 23, 2012

1 INTRODUCTION

Transmembrane β -barrel proteins (TMBBs) are important proteins that cross the lipid bilayer with a series of β -strands arranged in a cylindrical geometry and forming a structure that resembles a barrel (Schulz, 2000). Although all living organisms have transmembrane proteins organized into all- α helical bundles, TMBBs are presently found in the outer membranes of Gram-negative bacteria, mitochondria and chloroplasts. Despite their relevance, few TMBB structures are available at atomic resolution from Gram-negative organisms [about 0.1% of all the structures from Gram-negative organisms in the protein database (PDB)].

Several computational methods have been developed to predict TMBBs from protein sequences. Two prediction problems can be addressed: (i) TMBB detection in genomes (Bigelow *et al.*, 2004; Casadio *et al.*, 2003; Freeman and Wimley, 2010; Gromiha

and Suwa, 2006; Hayat and Elofsson, 2012; Remmert *et al.*, 2009; Savojardo *et al.*, 2011), and (ii) the prediction of protein topology (Bagos *et al.*, 2004, 2005; Bigelow *et al.*, 2004; Fariselli *et al.*, 2009; Hayat and Elofsson, 2012; Martelli *et al.*, 2002).

Here, we present BETAWARE, a tool based on machine-learning approaches to detect TMBBs in large sets of proteins and predict their topology. BETAWARE is available under GPL license as a standalone program, which is particularly well-suited for large-scale genome analyses. For the prediction of TMBB topology, it has been shown on comparative analyses that approaches based on grammatical modelling are the best performing (Fariselli *et al.*, 2009; Hayat and Elofsson, 2012), while for the detection of TMBBs in a set of proteins, the best-performing method is based on N-to-1 Extreme Learning Machines (ELM) (Huang *et al.*, 2006; Mooney *et al.*, 2011; Savojardo *et al.*, 2011).

2 METHODS

Detecting TMBBs in large sets of proteins is like finding a needle in a haystack. To address this problem, we used the previously developed method based on N-to-1 network encoding (Mooney *et al.*, 2011) and ELM training algorithm (Huang *et al.*, 2006). This method is the best performing on this task as compared with the most recent approaches (Savojardo *et al.*, 2011 and Supplementary Table S1). Once a protein sequence is predicted as putative TMBB, we applied a Grammatical Restrained Hidden Conditional Random Field (GRHCRF) model to predict the protein topology (Fariselli *et al.*, 2009). The model used is the same as previously described (Fariselli *et al.*, 2009), but here it has been retrained on a larger set of proteins.

2.1 Usage and program requirements

BETAWARE is written in pure Python to allow high compatibility. However, it has been explicitly designed to run on Unix/Linux systems (although it would be not too difficult to modify it for other operating systems). The BETAWARE program requires that the following packages are installed in the system: (i) python v2.x (tested on 2.6 or later), python argparse library, python numpy and scipy libraries (under Linux debian/ubuntu you have just to type on a single line: `sudo apt-get install python-numpy python-scipy python-argparse`). Once the software is downloaded and uncompressed, it is possible to run it directly moving into the BETAWARE root directory and typing the command:

```
$> ./predBeta.sh FASTA_FILE PROFILE_FILE
```

where `predBeta.sh` is a simple bash script, which takes as arguments a file containing the sequence to predict in FASTA format and its corresponding sequence profile in a separate file. An example of a

*To whom correspondence should be addressed.

