

# Empirical limits for template-based protein structure prediction: the CASP5 example

B. Contreras-Moreira<sup>\*,1</sup>, I. Ezkurdia, M.L. Tress, A. Valencia<sup>\*</sup>

*Protein Design Group, Centro Nacional de Biotecnología, Madrid, Spain*

Received 29 November 2004; revised 17 December 2004; accepted 5 January 2005

Available online 19 January 2005

Edited by Robert B. Russell

**Abstract** Most protein structure prediction methods use templates to assist in the construction of protein models. In this paper, we analyse the current state of template-based modelling approaches and reach an estimate of the empirical limits of these methods. Our analysis shows that current prediction methods are already reaching these empirical accuracy limits in the easier cases, where finding a close homologue to the native target structure is not a problem. However, we find that even in the absence of alignment errors and using optimal templates, template-based methods have intrinsic limitations, suggesting that other methodologies, such as *ab initio* procedures, must be used if accuracy is ultimately to be improved.

© 2005 Federation of European Biochemical Societies. Published by Elsevier B.V. All rights reserved.

**Keywords:** Template-based protein structure prediction; Comparative modelling; Fold-recognition; Fragment reconstruction; Accuracy limit

## 1. Introduction

Methods for protein structure prediction can be classified into two basic classes: those which use physical principles to fold a protein and those which use experimentally determined structures to help reconstruct the protein of interest. The first class is usually known as *ab initio* approaches [1]; the second includes related techniques such as comparative modelling, fold recognition and threading [2–7]. These generally use sequence alignments to map the sequence to be modelled onto protein templates of known structure and are guided by criteria such as sequence similarity or secondary-structure compatibility.

This paper deals mainly with the second class of methods, template-based methods. The empirical basis for these approaches comes from the observation by Chothia and Lesk [8] that protein sequence identity and structural similarity are correlated. According to their original results there are clear empirical limits for protein structure predictions based on sin-

gle templates: for proteins sequences around 95% identical backbone deviations are expected to be under 1 Å RMS; when the sequence identity drops to 30%, deviations grow to around 4 Å RMS. These limits broadly agree with the observed performance of comparative modelling servers as measured by continuous benchmarks such as EVA [9] (see [10] for a review), and ultimately affect the quality and therefore the applicability of template-based predictions [11].

In addition to these natural restrictions, methods for template-based prediction of protein structure must solve two technical problems: the choice of the template closer to the target structure, and the derivation of the sequence alignment between the query and template protein closer to the optimal structural alignment. The lack of satisfactory solutions for these two problems has been identified as negatively affecting the performance of fold recognition and comparative modelling methods in previous “Critical Assessment of Techniques for Protein Structure Prediction” experiments (CASP [12]) [13,14].

However, choosing the correct template and alignment are not the only problems facing predictors. Even those models built from the correct template and alignment often require substantial refinement in order to be sufficiently close to the native target structure. This paper seeks to estimate the limits of current template-based structure prediction techniques under ideal conditions, that is building a model *a posteriori* using multiple optimal templates and in the absence of alignment errors.

We do that by allowing models to be built by combining aligned fragments from several templates, selected by structural similarity. We then measure, using the CASP GDT\_TS score [15], how the best fragment-based predictions compare to the native target structure.

Additionally, we ask how far the predictions are from these best possible models. This gives us a better idea of how successful the current modelling methods are, how good they could be in the absence of the sequence alignment problem, and can implicitly tell us to what extent *ab initio* methods would be needed to improve the current performance of template-based methods.

## 2. Datasets, methods and algorithms

A collection of 68 targets, as split in domains by the CASP5 organisers, was taken as our test set (see <http://predictioncenter.llnl.gov/casp5>). These targets are proteins whose experimental structures were about to be released at the time CASP5 started (May, 2002). To model

<sup>\*</sup>Corresponding authors. Fax: +34 91 585 45 06.  
E-mail addresses: [contrera@cag.unam.mx](mailto:contrera@cag.unam.mx) (B. Contreras-Moreira),  
[valencia@cnb.unam.es](mailto:valencia@cnb.unam.es) (A. Valencia).

<sup>1</sup> Present address: Programa de Genómica Computacional, Centro de Ciencias Genómicas, Av. Universidad s/n, Colonia Chamilpa, 62210 Cuernavaca, Morelos, México.



Download English Version:

<https://daneshyari.com/en/article/10873685>

Download Persian Version:

<https://daneshyari.com/article/10873685>

[Daneshyari.com](https://daneshyari.com)