

MoDPepInt: An interactive webserver for prediction of modular domain-peptide interactions

Kousik Kundu^{1,†}, Martin Mann^{1,†}, Fabrizio Costa¹ and Rolf Backofen^{1,2,3,4 *}

¹Bioinformatics Group, Department of Computer Science, University of Freiburg, Germany, ²Centre for Biological Signalling Studies (BIOSS), University of Freiburg, Germany, ³Centre for Biological Systems Analysis (ZBSA), University of Freiburg, Germany, ⁴Center for non-coding RNA in Technology and Health, University of Copenhagen, Grønnegårdsvej 3, DK-1870 Frederiksberg C, Denmark

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: XXXXXXXX

ABSTRACT

Summary: MoDPepInt (Modular Domain Peptide Interaction) is a new, easy-to-use webserver for the prediction of binding partners for modular protein domains. Currently we offer models for SH2, SH3 and PDZ domains via the tools SH2PepInt, SH3PepInt and PDZPepInt. More specifically our server offers predictions for 51 SH2 human domains and 69 SH3 human domains via single domain models, and predictions for 226 PDZ domains across several species, via 43 multi-domain models. All models are based on support vector machines with different kernel functions ranging from polynomial, to Gaussian, to advanced graph kernels. In this way we model non-linear interactions between amino acid residues. Results were validated on manually curated data sets achieving competitive performance against various state-of-the-art approaches.

Availability: The MoDPepInt server is available under the URL: <http://modpepint.informatik.uni-freiburg.de/>

Contact: backofen@informatik.uni-freiburg.de

1 INTRODUCTION

Protein-protein interactions are often mediated by modular protein domains in eukaryotes and play an essential role in diverse biological processes such as signal transduction, cellular growth, cell polarity etc. (Pawson and Nash, 2003). Modular domains that specifically bind with short linear peptides are known as peptide recognition modules (PRMs). Each domain family recognizes peptides with specific characteristics. For example, phosphotyrosine (pY) containing peptides, proline-rich peptides and C-terminus peptides are recognized by SH2, SH3 and PDZ domains, respectively. However, individual domains from the same family show different binding specificity. Accurate models that can help understand the mechanisms responsible for the highly selective binding affinity are therefore of interest. Recently, several high-throughput techniques, such as protein microarray, phage display and SPOT synthesis, have been developed which can detect the binding specificity of various modular domains. However efficient

bioinformatics tools are needed in order to extract meaningful knowledge from the enormous amount of data produced.

To this end, we used state-of-the-art machine learning approaches to build support vector machine (SVM) models that can accurately predict binding specificity. We have collected into a unified web-based system called MoDPepInt, three different tools: SH2PepInt, SH3PepInt and PDZPepInt for three different modular domains, namely SH2, SH3 and PDZ (Kundu *et al.*, 2013b,a; Kundu and Backofen, 2014). Currently we offer single domain models for 51 SH2 human and 69 SH3 human domains, and multi-domain models for 226 PDZ domains across the species. To assess the quality of our models we have used manually curated interaction data achieving competitive performance against various state-of-the-art approaches.

In summary, MoDPepInt unique features include (i) the largest number of modeled domains, and (ii) a comprehensive SH2, SH3 and PDZ domain-peptide prediction system in a single platform.

2 APPLICATION AND FUNCTIONALITY

2.1 Input

All tools have a unified input format. Query sequences (up to a maximum number of 500) can be supplied either in a FASTA format or using UniProt database accession numbers. PDZPepInt offers predictions also for domains that are newly developed and/or not comprised in the original 226 PDZ domains: the unknown query domain should be supplied in FASTA format. Multiple query domain sequences can also be provided.

2.2 Filters

Several filters are available to increase predictive accuracy. SH2 domains generally recognize phosphotyrosine (pY) residues of binding proteins. For this reason in SH2PepInt we offer a *phosphotyrosine* filter that only considers those peptides whose tyrosine phosphorylation has already been experimentally verified and reported in PhosphoSitePlus database (Hornbeck *et al.*, 2012).

As SH3 domains mainly bind with proline rich peptides, in SH3PepInt, we offer a *proline rich* filter that uses 31 regular expressions to select proline rich peptides (Carducci *et al.*, 2012).

*to whom correspondence should be addressed.

† The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First Authors.

