

Aberystwyth University

Presaging critical residues in protein interfaces-web server (PCRPI-W)

Segura Mora, Joan; Assi, Salam A; Fernandez-Fuentes, Narcis

Published in:
PLoS ONE

DOI:
[10.1371/journal.pone.0012352](https://doi.org/10.1371/journal.pone.0012352)

Publication date:
2010

Citation for published version (APA):

Segura Mora, J., Assi, S. A., & Fernandez-Fuentes, N. (2010). Presaging critical residues in protein interfaces-web server (PCRPI-W): a web server to chart hot spots in protein interfaces. *PLoS ONE*, 5(8), e12352. <https://doi.org/10.1371/journal.pone.0012352>

General rights

Copyright and moral rights for the publications made accessible in the Aberystwyth Research Portal (the Institutional Repository) are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the Aberystwyth Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the Aberystwyth Research Portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

tel: +44 1970 62 2400
email: is@aber.ac.uk

Presaging Critical Residues in Protein interfaces-*Web Server* (PCRPI-*W*): A Web Server to Chart *Hot Spots* in Protein Interfaces

Joan Segura Mora¹, Salam A. Assi¹, Narcis Fernandez-Fuentes*

Section of Experimental Therapeutics, Leeds Institute of Molecular Medicine, University of Leeds, Leeds, United Kingdom

Abstract

Background: It is well established that only a portion of residues that mediate protein-protein interactions (PPIs), the so-called *hot spot*, contributes the most to the total binding energy, and thus its identification is an important and relevant question that has clear applications in drug discovery and protein design. The experimental identification of hot spots is however a lengthy and costly process, and thus there is an interest in computational tools that can complement and guide experimental efforts.

Principal Findings: Here, we present Presaging Critical Residues in Protein interfaces-*Web server* (<http://www.bioinsilico.org/PCRPI>), a web server that implements a recently described and highly accurate computational tool designed to predict critical residues in protein interfaces: PCRPI. PCRPI depends on the integration of structural, energetic, and evolutionary-based measures by using Bayesian Networks (BNs).

Conclusions: PCRPI-*W* has been designed to provide an easy and convenient access to the broad scientific community. Predictions are readily available for download or presented in a web page that includes among other information links to relevant files, sequence information, and a Jmol applet to visualize and analyze the predictions in the context of the protein structure.

Citation: Segura Mora J, Assi SA, Fernandez-Fuentes N (2010) Presaging Critical Residues in Protein interfaces-*Web Server* (PCRPI-*W*): A Web Server to Chart *Hot Spots* in Protein Interfaces. PLoS ONE 5(8): e12352. doi:10.1371/journal.pone.0012352

Editor: Ashley M. Buckle, Monash University, Australia

Received: April 21, 2010; **Accepted:** July 28, 2010; **Published:** August 23, 2010

Copyright: © 2010 Segura Mora et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Funding: This work was supported by the Research Councils United Kingdom Academic Fellow scheme (to NFF), Wellcome VIP award (to SAS), and an internal scholarship awarded by the Leeds Institute of Molecular Medicine (to JSM). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

* E-mail: N.Fernandez-Fuentes@leeds.ac.uk

† These authors contributed equally to this work.

‡ Current address: Faculty of Biological Sciences, Institute of Molecular and Cellular Biology, University of Leeds, Leeds, United Kingdom

Introduction

Cellular tasks require highly precise and regulated communication between proteins. Whether a protein is part of a metabolic pathway, an intermediate signalling effector, part of the transcription machinery, or a component of the cytoskeleton -just to mention some examples- requires proteins to act as complexes rather than as isolated units. Thus, protein-protein interactions (PPIs) are ubiquitous in Biology and therefore offer an enormous potential for the discovery of novel therapeutic agents able to modulate PPIs.

The analysis of protein complexes for which tertiary structure is known, has shown that protein interfaces are large, typically between 1500–2000 Å² [1,2], involving many intermolecular contacts (10 to 30 side chains per protein on average), and that such surfaces are usually flat and lacking defining physicochemical traits. It is for that reason that the identification of small-molecules that can act as modulators of PPIs is widely regarded as a formidable goal. However, as recently reviewed by Wells and McLendon [3] (and references therein), exciting new data

indicates that disruption of protein associations using small molecules is possible.

Part of the recent successes in the modulation of PPIs using small molecules has been possible by direct targeting of the important region, or *hot spot*, of the protein interface. The concept of hot spots in protein interfaces originates from the pioneering work of Clackson and Wells [4] that jointly with subsequent scientific works, have shown that most of binding energy in protein-protein associations can be ascribed to a small and complementary set of interfacial residues – a hot spot- surrounded by weaker interactions.

The experimental identification of hot spots in protein interfaces by Alanine scanning [5], Alanine shaving [6], or residue grafting [6], is a lengthy, labour-intensive, and costly process. Computational tools can be used to help and guide experimental efforts. We recently developed a novel computational tool: Presaging Critical Residues in Protein interfaces (PCRPI), that proved to be highly accurate and competitive with current computational methods [7]. In this paper, we present the implementation of the method as web application that will provide convenient and easy access to the

method to the scientific community. The web application has been designed having in mind a wide range of potential users, thus it has a user-friendly and straightforward interface with a minimal number of tunable parameters. Predictions are readily available for download or presented in a web page that has a number of functionalities such as a Jmol applet to visualize and analyze the predictions in the context of the protein structure.

Results and Discussion

Submitting a task

Running a prediction on PCRPI-*W* is a straightforward procedure. On the submission web page (Figure 1, panel A), users have to submit the coordinates of the protein complex of interest by either selecting it from a locally mirrored Protein Databank (PDB) database [8] typing the PDB code in a text box or uploading the coordinates (PDB format only); and select the chain identification code of the protein of interest. In advanced options, users can choose the type of BN and training set (see below).

Prior to prediction, structures undergo a set of quality checks. If atoms present alternative locations or rotamers, only the first occurring rotamer is kept. Also, if residues have insertion codes,

the distance with neighboring residues is calculated and discarded if structurally equivalent. Side-chains with missing atoms are reconstructed using Scrwl 4.0 [9], an important step because energy calculations are highly affected by missing atoms. Finally, the length of proteins are checked and those shorter than 40 residues are discarded. As a result, a modified version of the original coordinate file, remediated coordinates file, is generated. This is the file used as input during the prediction and is downloadable from the result web page. Changes to the original coordinate (if any) are recorded in the log file (see below *Retrieving and visualizing results*).

PCRPI-*W* features two types of BNs, a naïve and expert, that can be trained using two different datasets: Ab+ and Ab- (Figure 2). More information about the structure of the BNs and the composition of the training sets can be found in the help web page of the server or in the original publication describing the method [7]. By default, PCRPI-*W* run the prediction using a naïve version of the BN trained on the Ab+ dataset, although both, BNs type and training sets are tunable parameters and users can select the ones that adjust the best to their needs. If an e-mail address is given at time of the submission, user will be notified by e-mail once the job is finished including a hyperlink to the results web page

Panel A: Submission Page

Protein complex

PDB code: 1pky

OR

upload your own coordinates (PDB format) no file selected

Chain ID: A

Prediction parameters

Bayesian network architecture: naïve

Training dataset: Ab+

Enter your e-mail address (optional): sisux@yahoo.com

Retrieve your PCRPI prediction

Enter job id (as provided after submission):

Panel B: Confirmation Page

Your prediction has been submitted!

Your assigned job ID is: PCRPI_782070659

Here your prediction parameters:

PDB file: 1pky (PDB databank)

Selected chain: A

Bayesian network architecture: naïve

Panel C: Results Page

Prediction has finished successfully!

[Prediction parameters] [Download files] [Results]

Status for job id PCRPI_591751327 is: FINISHED

Prediction parameters:

Submitted protein complex (coordinates)

Selected chain: A

Bayesian network architecture: naïve

Training dataset: Ab+

Generated files:

Output coordinates (B-factor = prediction probability * 100)

Log file

Prediction file (plain text tab-delimited)

Predicted critical residues (sorted by probability):

Chain ID	Residue type	Residue number (as in PDB file)	Probability	Scaled probability (Min-Max normalization)	Scaled_prob(i)=(prob(i)-Min_prob)/(Max_prob-Min_prob)
E	192	1.0000	1.0000		
D	215	1.0000	1.0000		
D	45	0.9953	0.9953		
K	46	0.9944	0.9944		
E	203	0.9942	0.9942		
E	155	0.9885	0.9885		
W	58	0.8627	0.8623		
D	212	0.8345	0.8340		
H	211	0.7795	0.7788		
R	194	0.6712	0.6701		
A	196	0.6447	0.6436		
R	164	0.5842	0.5849		
L	193	0.4737	0.4730		
P	50	0.4704	0.4687		
P	219	0.4491	0.4474		
A	67	0.4523	0.4505		
V	104	0.4091	0.4072		
S	210	0.3458	0.3487		
P	48	0.1999	0.1973		
L	154	0.1875	0.1849		

Panel D: Jmol Applet

Protein complex. If you know the job ID you can use the job ID to mark the following URL:

Panel E: Jmol Applet

RESID	TYPE	PROB
192	E	1
215	D	1
45	D	0.9953
46	K	0.9944
203	E	0.9942
155	E	0.9885
58	W	0.8627
212	D	0.8345
211	H	0.7795
200	T	0.7548
194	R	0.6712

Figure 1. Several screenshots of PCRPI-*W*. The home web page of the server is the submission web page (A), where upon submission a temporary web page (B) reports an unique job identification code and a link to the results web page that users can bookmark to retrieve their results when available. The results web page (C) provides access to a number of links among them: a link to download the list of predicted hot spot residues (D) and a link to visualize the protein complex colored by prediction probabilities using a Jmol applet (E). doi:10.1371/journal.pone.0012352.g001

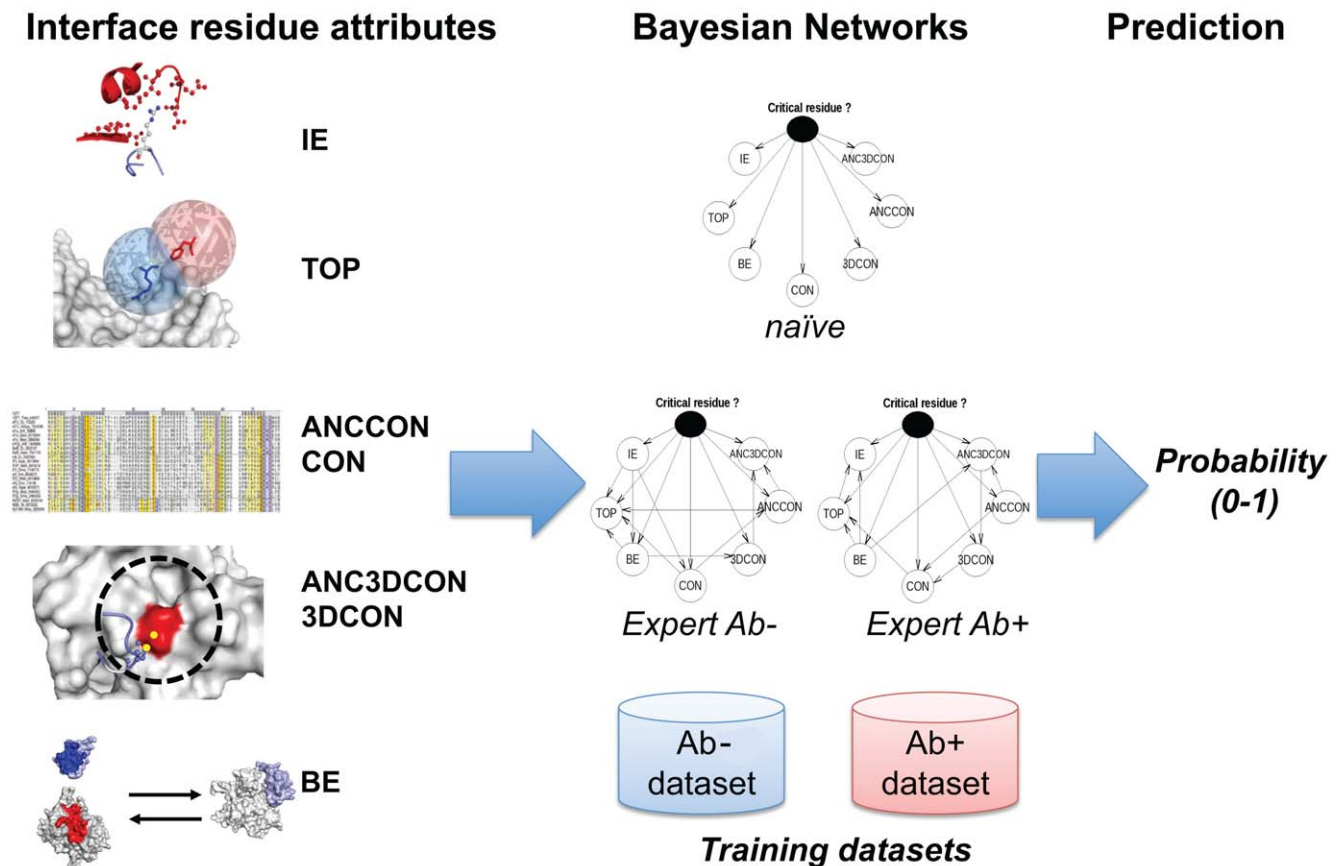


Figure 2. General overview of the prediction process. PCRPI combines seven different measures by using BNs and outputs a probability. The input variables are: IE, TOP, BE, CON, 3DCON, ANCCON, and ANC3DCON. There are two different training datasets: Ab+ and Ab−, and three different BNs: a naïve and two training dataset-specific experts BNs that can be invoked during the prediction. For more information regarding PCRPI method and input variables, refer to the original publication describing the method [7].
doi:10.1371/journal.pone.0012352.g002

(hyperlink also shown upon submission for bookmarking purposes; Figure 1, panel B). PCRPI-*W* assigns a unique job identifier for each submitted job (e.g. PCRPI_cA8r0nAz0). This job identifier can be used to check the status of the submission (i.e. in queue, running, finished) and to retrieve the results by typing it in the 'Job ID' field at the submission web page.

Jobs are handled by a queuing system and, if not competing jobs, typically take few minutes to be completed; larger protein complexes featuring large or multiple interfaces can take up to one hour. The most time consuming is the estimation of the binding free energy, which for large interfaces and protein complexes requires intensive and long computational times, and the sequence search and calculation of sequence profiles for evolutionary-based measures.

Retrieving and visualizing results

PCRPI-*W* returns a list of interface residues sorted by probability (Figure 1 panel C and D) and several links to download files used or generated during the prediction. A successful prediction will generate the following files: a file that contains the original coordinates as uploaded by the user or as in the PDB; the remediated version of the coordinates file (see above *submitting a task*); a modified version of the input coordinates where the B-factor field has been substituted by a value that is equal to the prediction probability times 100 (facilitating analysis of predictions when using molecular visualization programs such as

PyMOL [10]); a list of interface residues sorted by probability; a file detailing the atomic interaction of the interface residues as defined by CSU program [11] (atomic interactions can be also visualized in the context of the structure by using a Jmol (<http://www.jmol.org>) applet, see next); and a log file that records the entire prediction process and that can be examined if errors are reported.

Other elements that are shown in the results web page is the mapping of predictions on the protein sequence and a Jmol applet that allows the visualization of the structure of the complex and the mapping of the predictions. The Jmol applet includes a clickable list of protein chains and residues sorted by probability (Figure 1, panel E), and thus facilitate the process of visualization and selection of interface residues and predictions. Upon selection of a given residue, this will be highlighted in ball-and-stick representation and the atomic interactions with neighbouring residues will be shown.

Possible bottlenecks

Occasionally, PCRPI-*W* may fail to provide a prediction. The main reason is usually when the coordinates file contains only one protein chain or if more than one, these do not interact, i.e. no atomic interactions between protein chains. In this case, interface(s) cannot be located and therefore the program fails. More rarely, there can be errors along the prediction process, e.g. problems during free energy calculations or errors when deriving

Table 1. Comparison of different methods for the prediction of critical residues in protein interfaces using a BiD derived dataset as described in Tuncbag et al. [18].

Method	Precision (P)	Recall (R)	F1 score
PCRPI ^a	0.79	0.64	0.71
FoldX ^b	0.75	0.36	0.49
Robetta-Ala ^c	0.63	0.57	0.60
KFC ^c	0.51	0.36	0.42
KFC-A ^c	0.53	0.48	0.51
LDA ^c	0.72	0.57	0.64
Tuncbag et al. [18] ^c	0.73	0.59	0.65

^aPredictions were performed using PCRPI [7] with an expert BN trained in a Ab+ dataset that does not include the crystal structure of the c2 fragment of streptococcal protein G in complex with the Fc domain of human Ig (PDB code 1fcc).

^bValues were obtained running FoldX [19] with default parameters and a $\Delta G_{\text{binding}}$ cut-off of 2.0 Kcal.mol⁻¹ (i.e. residues were considered critical if upon mutation to Ala, predicted $\Delta G_{\text{binding}} \geq 2.0$ Kcal.mol⁻¹).

^cPrecision, recall, and F1 score values taken from Tuncbag et al. [18].
doi:10.1371/journal.pone.0012352.t001

evolutionary-based measures, e.g. PSI-BLAST [12] fails to find homologous sequences with significant E-values. As described above, a log file is available for users to download and examine to understand the reason(s) of reported error(s). In addition, users can contact the authors via e-mail for further support.

Availability and Future Directions

PCRPI-W server is freely available upon registration to the scientific community at <http://www.bioinsilico.org/PCRPI>. Besides the option of submitting tasks to the server, users can browse an extensive documentation, have access to related resources available online, and download the benchmark and training datasets.

Methods

Prediction algorithm

Several are the features that characterize the residues that are part of a hot spot and these have been exploited in the past for prediction purposes. These features can be broadly grouped in three

References

- Jones S, Thornton JM (1996) Principles of protein-protein interactions. *Proc Natl Acad Sci USA* 93: 13.
- Lo Conte L, Chothia C, Janin J (1999) The atomic structure of protein-protein recognition sites. *J Mol Biol* 285: 2177.
- Wells JA, McClendon CL (2007) Reaching for high-hanging fruit in drug discovery at protein-protein interfaces. *Nature* 450: 1001–1009.
- Clackson T, Wells JA (1995) A Hot Spot of Binding Energy in a Hormone-Receptor Interface. *Science* 267: 383–386.
- Wells JA (1991) Systematic mutational analyses of protein-protein interfaces. *Methods Enzymol* 202: 390–411.
- Jin L, Wells JA (1994) Dissecting the energetics of an antibody-antigen interface by alanine shaving and molecular grafting. *Protein Sci* 3: 2351–2357.
- Assi SA, Tanaka T, Rabbitts TH, Fernandez-Fuentes N (2009) PCRPI: Presaging Critical Residues in Protein interfaces, a new computational tool to chart hot spots in protein interfaces. *Nucleic Acids Res* 38(6): e86.
- Berman HM, Battistuz T, Bhat TN, Bluhm WF, Bourne PE, et al. (2002) The Protein Data Bank. *Acta Crystallogr D Biol Crystallogr* 58: 899–907.
- Wang Q, Canutescu AA, Dunbrack RL, Jr. (2008) SCWRL and MolIDE: computer programs for side-chain conformation prediction and homology modeling. *Nat Protoc* 3: 1832–1847.
- <http://www.pymol.org/> (last accessed 2010).
- Sobolev V, Sorokine A, Prilusky J, Abola EE, Edelman M (1999) Automated analysis of interatomic contacts in proteins. *Bioinformatics* 15: 327–332.
- Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, et al. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25: 3389.
- DeLano WL (2002) Unraveling hot spots in binding interfaces: progress and challenges. *Curr Opin Struct Biol* 12: 14–20.
- Pearl J (1988) *Probabilistic Reasoning in Intelligent Systems*. San Francisco: Morgan Kaufmann Publishers.
- Jordan M (1998) *Learning in Graphical Models*. London: The MIT Press.
- Tanaka T, Williams RL, Rabbitts TH (2007) Tumour prevention by a single antibody domain targeting the interaction of signal transduction proteins with RAS. *EMBO J* 26: 3250–3259.
- Pruitt KD, Tatusova T, Maglott DR (2007) NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res* 35: D61–65.
- Tuncbag N, Gurosoy A, Keskin O (2009) Identification of computational hot spots in protein interfaces: combining solvent accessibility and inter-residue potentials improves the accuracy. *Bioinformatics* 25: 1513–1520.
- Guerois R, Nielsen JE, Serrano L (2002) Predicting changes in the stability of proteins and protein complexes: a study of more than 1000 mutations. *J Mol Biol* 320: 369–387.

categories depending on nature of the data. Hot spots can be predicted by energy, structural, and evolutionary-based (e.g. sequence conservation) analysis. Although these features are useful, it was shown that, individually, cannot unambiguously define hot spots [13]. PCRPI [7] overcomes this limitation by combining a set of seven different measures that account for energetic, structural, and evolutionary-based information (Figure 2). Individual measures are combined into a unique probabilistic framework by using Bayesian Networks (BNs) [14,15].

The performance of PCRPI was benchmarked in two independent datasets [7]. The first set was composed of 25 protein complexes summing up 636 interfaces residues, 300 of which were validated as critical or non-critical residues by experimental means and available in the scientific literature. The second dataset was the protein complex formed by HRAS and a VH domain of an Fv antibody [16]. Under both scenarios PCRPI delivered highly accurate and consistent predictions. Moreover, in a head-to-head comparison with other available computational tools using the same test set, PCRPI predictions were superior in terms of precision, recall, and F1-scores (Table 1).

Design, implementation and use of PCRPI-W

PCRPI-W is implemented on an Apache server running on a Red Hat® enterprise linux-based operating system. The server is interfaced with a CGI Perl and Javascript coded web interface. PCRPI-W modules and accessory scripts are coded in Perl, Fortran, and C++ respectively. Databases required by the server, namely, PDB [8] and NCBI non-redundant (NR) protein sequence database [17], are locally mirrored and weekly updated. All the queries are submitted to a queuing system that submits the tasks to a computer farm. Results are displayed in HTML format and send to the user by e-mail containing a hyperlink to the results web page.

Acknowledgments

NFF thanks Dr. Gendra for critical reading and insightful comments to the manuscript and Ms Martina and Ms Daniela G Fernandez for continuing inspiration and motivation. Authors acknowledge constructive comments from anonymous reviewers.

Author Contributions

Conceived and designed the experiments: NFF. Performed the experiments: JSM SAA. Analyzed the data: NFF. Contributed reagents/materials/analysis tools: JSM SAA NFF. Wrote the paper: JSM SAA NFF.