Evaluation of currently available web-servers for 3D-structure based prediction of protein interface regions

Evan Guiney

Many important problems in molecular biology today involve characterizing protein interaction networks. There are numerous experimental approaches to this problem, ranging from narrow, focused biochemical analysis of direct physical interactions, to broader approaches like yeast two hybrid and mass spec; direct observation of an interface with a 3D crystal structure is the gold standard, but is also the most labor intensive approach. Though they are relatively new, computational approaches may also become an important additional tool.

Protein interfaces can be divided into two broad classes- obligate, strong interactions, such as usually found in homodimers, or between members of a stable complex, and transient interactions, for example between kinases and phosphatases and their substrates. The residues that make up the faces of an obligate interaction resemble the hydrophobic core of a globular protein, while residues involved in transient interactions have properties that can be thought of as intermediate between core and surface exposed regions. These unique characteristics of interface residues, as opposed to non-interface surface residues, make interface prediction a good computational problem, and there are currently several related methods for predicting interface regions (Zhou and Qin 2007, Neuvirth *et al.* 2004, Headd *et al.* 2007, Via *et al.* 2000, Yan *et al.* 2008, Yan *et al.* 2004).

In this paper, I will review 7 freely available, web-based servers for interface prediction, from the perspective of a molecular biologist. I will cover the general principles behind interface prediction, give a non-technical treatment of the statistical prediction methods, and review the pros and cons of each method. I will also evaluate the methods' performance predicting interfaces on PP2B/calcineurin (CN). CN is a challenging problem; it is an obligate heterodimer with two biochemically characterized interface regions, and was not included in any of these methods' training sets; our lab is currently investigating both known interface regions and looking for potential new regions.

The urls for the programs reviewed here are as follows:

Meta-PPISP    http://pipe.scs.fsu.edu/meta-ppisp.html (Qin and Zhou 2007b)

Cons-PPISP    http://pipe.scs.fsu.edu/ppisp.html (Chen and Zhou 2005)

ProMate       http://bioinfo.weizmann.ac.il/promate/ (Neuvirth *et al* 2004)

PPI-PRED      http://bmbpcu36.leeds.ac.uk/ppi_pred/ (Bradford and Westhead, 2005)

SPPIDER       http://sppider.cchmc.org/ (Porollo and Meller, 2007)

InterProSurf  http://curie.utmb.edu/pdbcomplex.html (Negi *et al.* 2007)

WHISCY        http://nmr.chem.uu.nl/Software/whiscy/index.html (de Vries *et al.* 2006)

All currently available prediction tools use some set of the following characteristics of interface residues:

- Conservation: Interface residues are highly conserved, while non-interface surface residues are under less selective pressure and are free to vary. Sequence conservation is used in every interface prediction program reviewed here.

- Residue bias: Although transient and obligate interfaces differ, which can present complications, hydrophobic residues are favored in interfaces, and charged residues are disfavored, except for arg, which often forms cation-$\pi$ interactions. (Headd *et al.* 2007, Zhou *et al.* 2007, Neuvirth *et al*, 2004)

- Solvent accessibility: Interface residues generally have higher solvent accessibilities than non-interface residues. Regions that never form interfaces have been selected to minimize unfavorable solvent interactions, while interface residues must balance solvent interactions in the unbound state with intermolecular interactions formed upon binding (Zhou *et al.* 2007, Jones and Thorton 1997).

- Entropy: Interface residues have lower temperature factors than other surface residues (Zhou *et al.* 2007, Neuvirth *et al.* 2004); this can be understood in terms of entropy cost upon binding: because interface residues have lower entropy to begin with, this minimizes the entropy cost incurred upon binding.

- Patches: Interfaces occur in roughly circular patches, so all evaluated prediction programs take into account the characteristics of surrounding residues.


Interface prediction methods:

The interface prediction tools reviewed here all require a solved 3D structure of the protein of interest (although most allow uploading of custom PDB files); interface

prediction methods based solely on primary sequence are available, but will not be discussed here. Generally speaking, the prediction programs score surface residues for some or all of the above properties, and find clusters of high scoring residues, which are identified as the interaction region.

The methods for weighing the different residue characteristics vary: Two programs (Whiscy and InterProSurf) use a relatively simply linear regression method. This has the advantage of being simple and transparent; Whiscy allows the user to choose which characteristics to include (conservation, residue bias and contiguity are the options), but gives the poorest performance.

Three methods use machine learning approaches- cons-PPISP and SPPIDER use neural networks and PPI-Pred uses Support Vector Machines. These methods are robust, but non transparent. Both utilize a training set of proteins with known interface and non-interface residues and seek to maximize correct predictions. Neural networks consist of input, 'hidden' and output nodes, linked by weighted linear transformations; residue characteristics (again, conservation etc) are inputs, and the weights between the nodes are adjusted to give optimal correct predictions. Similarly, the Support Vector Machine method maps residue characteristics to a high-dimensional space and finds a hyper-plane that best separates that space between interacting and non-interacting residues (Zhou and Qin 2007).

ProMate uses a naïve Baysian approach to evaluate, for each residue characteristic, what the probability of being in an interface is. This allows ProMate to be quite transparent; the user can select which characteristics to include in the analysis, from the following impressive list: Single amino acids distribution, atoms distribution,

chemical character, amino acid pairs distribution, evolutionary conserved positions, non-regular secondary structure length, sequence distances within a circle, secondary structure, domains, hydrophobic patch size/rank, temperature factor, or water molecules. This versatility is admirable, but without extensive knowledge of each term, it is hard to justify inclusion/exclusion and so I suspect most users will simply use the default selection.

Finally, meta-PPISP combines scores from cons-PPISP, PROMATE, and PINUP (PINUP is no longer supported on the web, and uses an empirically weighted combination of residue entropy, solvent accessibility, and conservation (Liang *et al 2006)*. The raw code (WARNING: not an .exe file) can be downloaded at http://sparks.informatics.iupui.edu/index.php?pageLoc=Publications).

Interface and data output evaluation:

All of these servers are very user friendly. Since prediction is based on 3D structure, all you have to do, in most cases, is enter your email and desired PDB id. Cons-PPISP and meta-PPISP allow you to evaluate multiple chains at once; all of the other servers can only handle a single chain. This is perhaps the most important functional difference between the programs, at least for multimeric proteins like calcineurin: while the PPISP servers can find surface interface regions, most of the other programs simply locate the interface between the A and B subunits, something anyone could do simply by looking at the PDB!

Another advantage of the PPISP servers is that they output clusters of residues, with confidence scores for the entire cluster, in addition to giving a interface/non-

interface score for each residue. PROMATE and Whiscy both also output scores for individual residues, but do not group clusters; PROMATE outputs as altered temperature factors in a PDB, making qualitative appraisal easy, while Whiscy gives a list. Interprosurf gives a list of predicted interface residues, and attempts to evaluate the change in exposed surface area for each residue upon binding, but this seems excessively specific and less useful that the more general scores given by the other prediction programs. PPI-PRED does not give individual residue scores, which is unfortunate, but does attempt to cluster residues (unfortunately, for CN, these clusters covered ~80% of the surface!), while SPPIDER simply outputs a PDB with residues marked as either interface or non-interface, with no clustering or individual scores.

SPPIDER and the PPISP servers both have good email notification of results, which allows easy access to past results; the other servers display results on a webpage. Most servers were quite fast, predictions for CN all took five minutes or less, except for meta-PPISP which took 4 days. Cons-PPISP can be accessed in two ways: either directly, at the url listed above, or as bundled in meta-PPISP. A meta-PPISP search will quickly return its cons-PPISP component, and then after a longer wait (again, for me, it was 4 days), return the full meta analysis. I found that the cons-PPISP results obtained from the meta-PPISP differed slightly from those from the cons-PPISP server; the cons-PPISP results from the meta-PPISP server were better, and are the results discussed here. These interface observations are ranked and summarized in table 1.

Prediction quality:

Zhou and Qin 2007 is an excellent review of cons-PPISP, meta-PPISP, ProMate, PPI-Pred, SPPIDER and the now-unsupported PINUP. Zhou and Qin summarize each program and then compare Coverage (True Positives/Total Interface Residues) and Accuracy (True Positives/True Positives + False Positives) for two large protein interaction datasets; Enz35, a set of 35 enzyme/inhibitor compelexs, and CAPRI, a 25 protein dataset from the protein interface prediction version of CASP. They find that for both the easy Enz35 and hard CAPRI datasets, the programs perform with a consistent ranking of, best to worst, meta-PPISP > PINUP > ProMate & cons-PPISP > SPIDDER > PPI-Pred (Zhou and Qin are also the authors of the PPISP programs). The ranking data is reproduced below, in figure 1, and summarized in table 1.

I have also conducted a focused evaluation of each prediction program for calcineurin, the protein our group studies. Calcineurin is an excellent candidate for such an evaluation. While much cell signaling takes place through cascades of kinases, which usually have well defined phosphorylation motifs, making substrate prediction possible from primary sequence alone, phosphatases also play key signaling roles but have more challenging substrate interactions. Phosphatases like calcineurin do not recognize a consensus dephosphorylation site; instead, specificity is achieved by docking interactions. Our lab has characterized two docking motifs used by yeast calcineurin, and shared with human calcineurin (for which several structures have been solved): PxIxIT (thought to be a primary docking site) and YLxVP (binding here is suspected to allosterically activate calcineurin).

Calcineurin is composed of a catalytic subunit, CNA, and a regulatory subunit, CNB. CNA also has a long, unstructured C-terminal regulatory tail, containing an

autoinhibitory domain that binds to the active site. Calcineurin is activated by calcium; calcium bound calmodulin binds to the a calmodulin binding domain in the regulatory tail (this domain is unstructured in all CN crystal structures), this removes an auto-inhibitory domain. Furthermore, calcium binding to CNB causes final, full activation of calcineurin. Docking sites have been mapped for YLxVP by mutagenesis and for PxIxIT by both mutagenesis and crystal structure. Furthermore, the inhibitor complexes FK506/FKBP and Cyclosporin/Cyclophilin bind calcineurin near the YLxVP binding region (Figure 2).

The different webservers had varying degrees of success with finding calcineurin binding sites. Two servers (cons-PPISP and meta-PPISP) were able to evaluate the entire calcineurin heterodimer, but the rest can only handle a single chain; in these cases analysis CN-A is most fruitful.

**Cons-PPISP** was the best performing server (Figure 4), its top scoring cluster included 5/6 verified PxIxIT binding residues, and its second hit included 3/3 verified YLxVP binding residues. Furthermore, because the PPISP servers can handle multimeric proteins, cons-PPISP did not return the most obvious, uninteresting interface, between CN-A and CN-B. Cons-PPISP also found potential interfaces on CN-B, as well as interesting sites of potential new docking regions, which will be discussed later (figure 10).

**Meta-PPISP** gave the next best results, probably simply due to its cons-PPISP component (Figure 3). Meta-PPISP found both the PxIxIT and YLxVP interfaces, although with low scores, but gave a high score to an interesting new region that was also identified by ProMate and SPPIDER (see figure 10).

**ProMate** found two potential docking regions; one was similar to the region identified by meta-PPISP and SPPIDER. Unfortunately, because ProMate can only handle a single protein chain, its other potential docking site is probably incorrect, as it is occluded by the CN-A CN-B interaction. However, ProMate did not predict that the CN-A alpha helix formed a docking site for CN-B.

**SPPIDER** performed marginally better than ProMate, but worse than the PPISP servers. It found 1/6 PxIxIT residues, and found the entire CN-A alpha helix that forms the CN-B binding region. SPPIDER also found the potential new binding site (see figure 10).

**InterProSurf** had less success, it was only able to identify the CN-A/CN-B interacting region, and did not find any other interaction patches on CN-A.

Finally, **Whiscy** and **PPI-pred** did not give acceptable results. PPI-pred was too permissive- it predicted interaction regions that cover nearly the entire surface of CN-A. Whiscey was less permissive, but it only found residues in the CN-A/CN-B interface, and the residues it identified there were only a small subset of CN-A/CN-B interacting residues.

New interface regions:

Three programs predicted the same general area of calcineurin as being a new protein-interaction site (orange patch in Figure 10). This region is not a known docking site, but has several interesting properties. Our lab has done a small amount of directed mutagenesis here, but not on any of the key residues identified by interface prediction. Our attention was first drawn to the region because it is well conserved, including in

Protein Phosphatase 1, the nearest homolog of calcineurin; in PP1 this region is indeed a substrate docking region, for the motif MyPhoNE (Roy and Cyert 2009). The yeast homolog of one residue identified by the interface predictors, D229, was mutated in a previous mutagenesis screen (unpublished results), and was found to drastically destabilize calcineurin.

Because three prediction programs (SPPIDER, PROMATE and meta-PPISP) implicate this region, further study in the lab is warranted. Y224, H250 and C256 were each identified by 2/3 programs; we are now in the process of designing primers to mutate these three residues to alanine, and will look for any effects on calcineurin phenotypes.

Cons-PPISP, which had the best results overall, found two addition potential interaction regions nearby (pink and yellow in Figure 10), and these regions had similar scores to the PxIxIT and LxVP regions.

Summary:

Table 1 summarizes the features of the interface prediction webservers reviewed here. Cons-PPISP performed best for evaluating calcineurin, while in general meta-PPISP has an edge. Because searches on the meta-PPISP server return, first, a cons-PPISP output, and then later, a combined meta-PPISP output, the meta-PPISP server is my number one choice for protein interface prediction. PROMATE also performs well, has the most flexibility in which features are evaluated, and has a good, quantitative output.
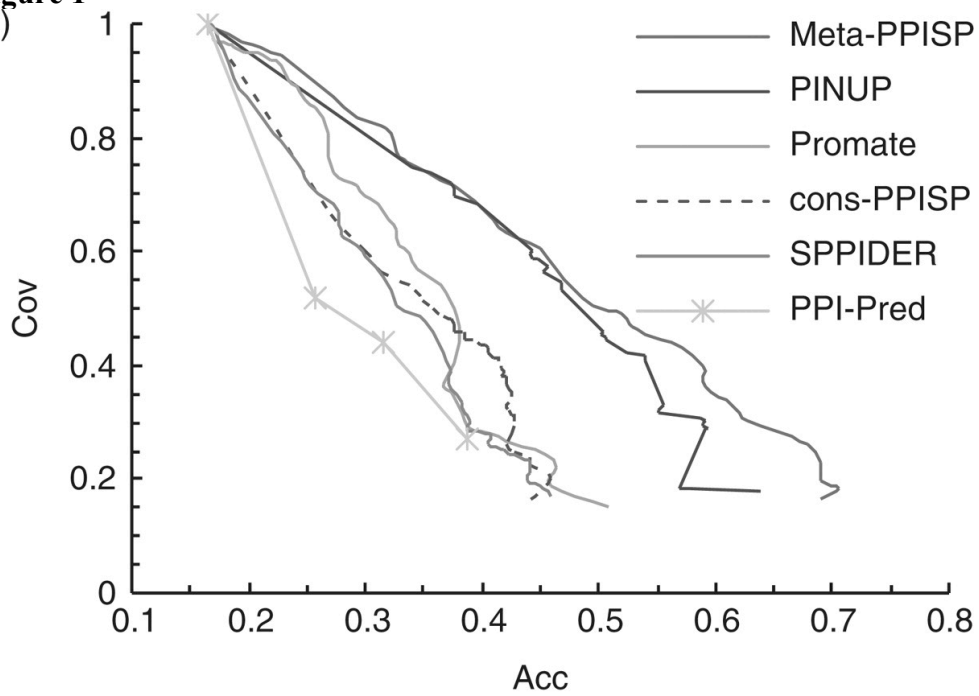
SPPIDER performed decently for calcineurin, but its output (residues are not grouped into clusters, and are categorized simply as interfact/not interface) is not ideal.

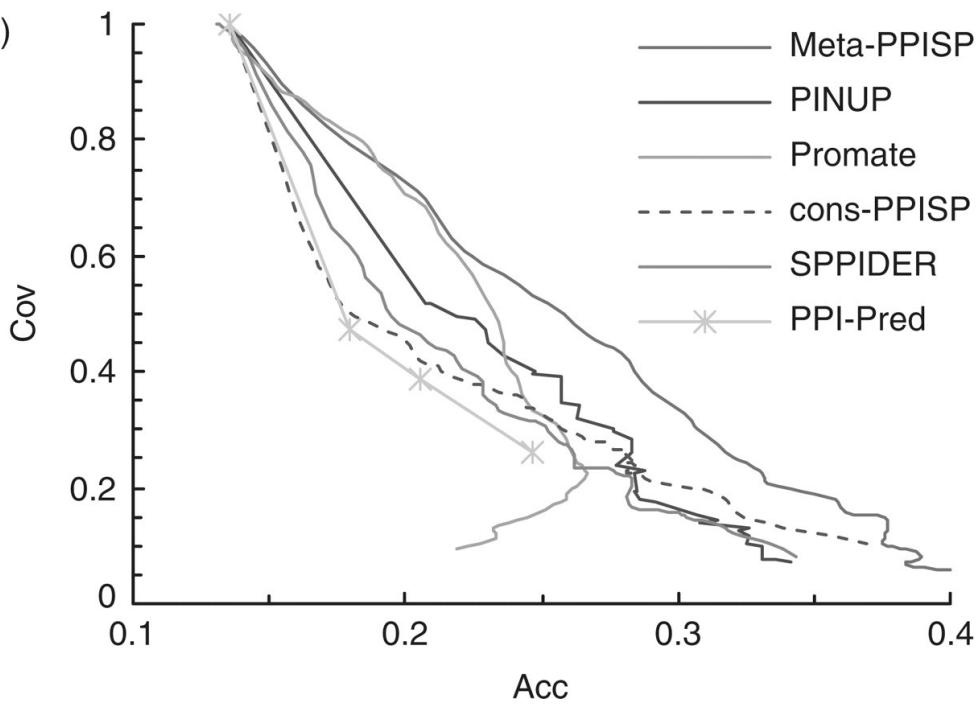PPI-Pred, Whiscy and InterProSurf are not recommended.

Table 1: Comparison of 3D-sturcture based interface prediction Webservers

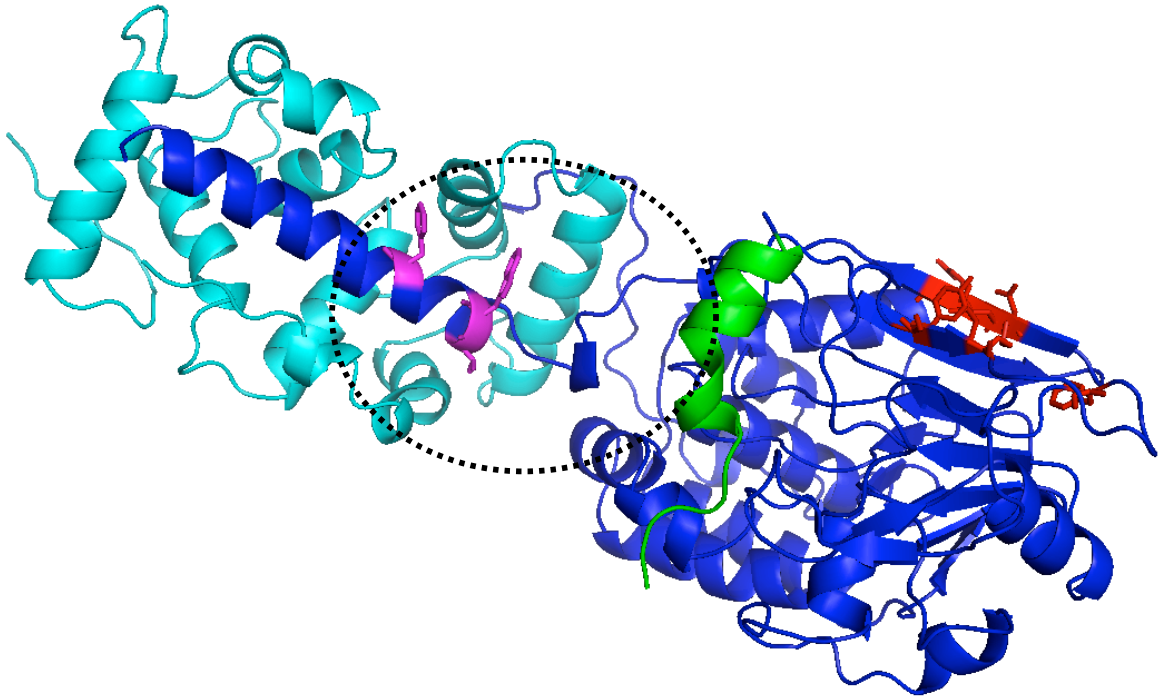| | Coverage/ Accuracy Rank (Zhou and Qin 2007) | Calcineurin docking sites | Multimer/ monomer only? | Output format | Other: PROs | Other: CONs |
|---|---|---|---|---|---|---|
| Meta-PPISP | 1 | Good (+PxIxIT, +YLxVP, unverified#1) | Multimers and monomers | Quantitative list, PDB-temperatures | Email output, gives both meta- and cons-PPISP results | Slow, no custom interface properties |
| Cons-PPISP | 3 | Best (++PxIxIT, ++YLxVP, unverified#2,3) | Multimers and monomers | Quantitative list, scored clusters | Email output | No custom interface properties |
| ProMate | 3 | Poor (unverified #1) | Single chain only | Quantitative PDB-temp | Highly customizable interface properties | Web output only |
| SPPIDER | 4 | Fair (+/-PxIxIT, CN-A/CN-B, unverified#1) | Single chain only | Binary PDB-temp | | Web output only, no custom interface properties |
| InterProSurf | - | Poor (CN-A/CN-B) | Single chain only | Semi-quantitative list | | Web output only, no custom interface properties |
| PPI-Pred | 5 | Failed | Single chain only | Binary PDB-temp, list | | Web output only, no custom interface properties |
| Whiscy | - | Failed | Single chain only | Quantitative list | Some interface properties customization, Use your own multiple alignments | Web output only |

**Figure 1**



(a) performance of webservers on Enz35 set
(b) performance of webservers on CAPRI set
From Zhou and Qin 2007

**Figure 2**
Caclineurin Sturcture

Calcineurin structure:
Calcineurin A subunit is in dark blue, B subunit is in light blue. The CN-A C-terminal autoinhibitory domain, is in green, and is connected to the rest of CN-A by an unstructured linker.
Residues verified *in vivo* to be involved in PxIxIT motif binding are in red; residues verified *in vivo* to be involved in YLxVP binding are in purple.

**Figure 3**
Meta-PPISP

AID

YLxVP 3/3
(low scores)

PxIxIT 6/6
(Intermediate scores)

**Figure 4**
cons-PPISP



PxIxIT: 5/6
verified *in vivo*

YLxVP: 3/3
verified *in vivo*

cons-PPISP clusters 1 (red) and 2 (green) (confid=17)



90°

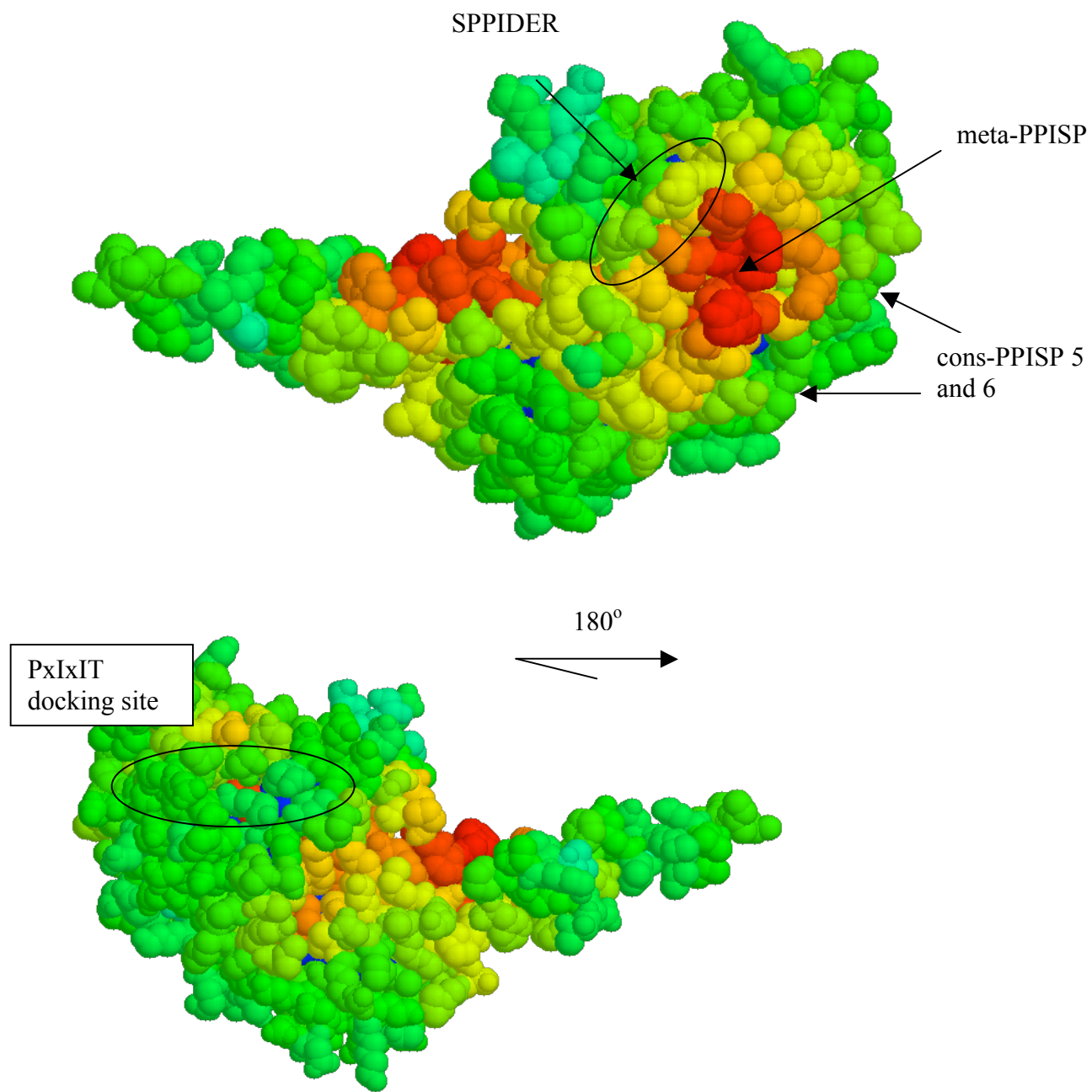cons-PPISP clusters 3 (yellow), 4 (pink), 5 (orange) and 6 (white). Confid=17,17,16,16
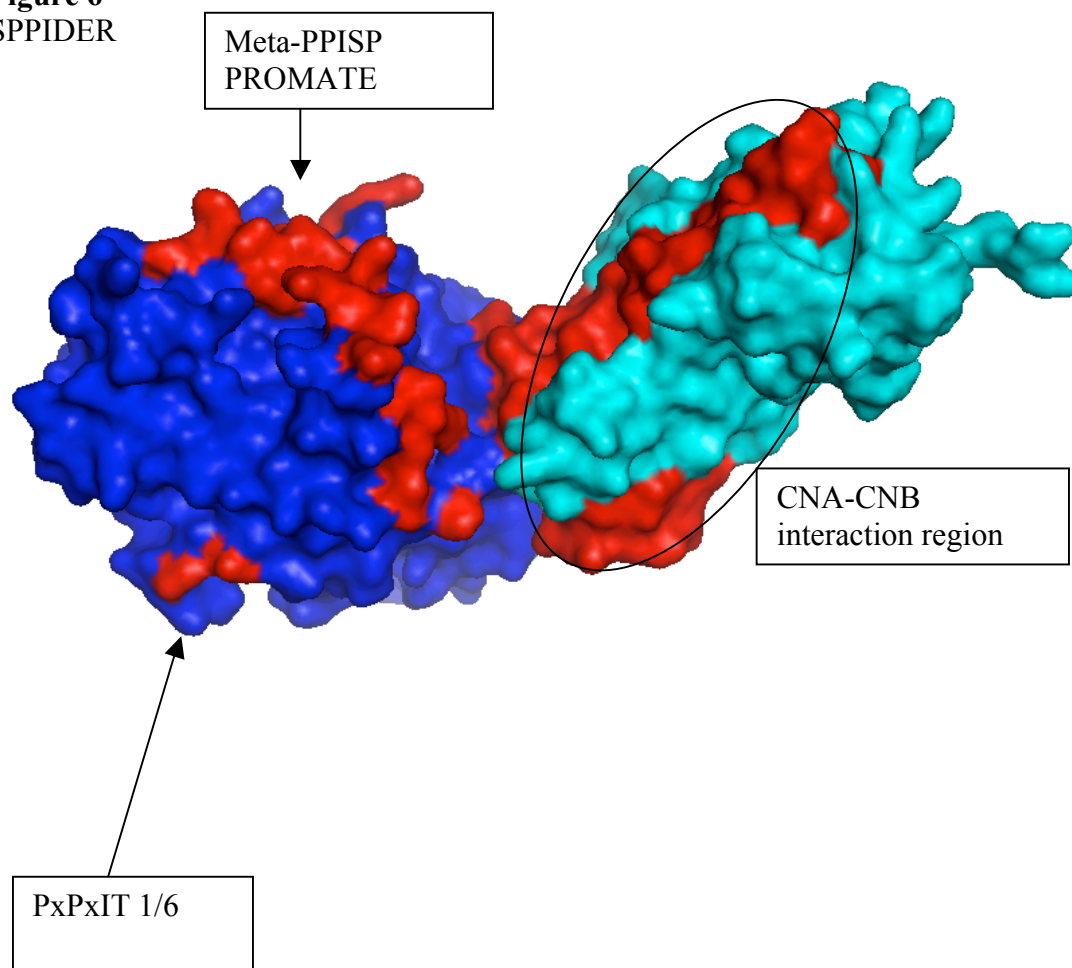
**Figure 5**
Promate

SPPIDER

meta-PPISP

cons-PPISP 5
and 6

180°

PxIxIT
docking site

**Figure 6**
SPPIDER

Meta-PPISP
PROMATE

CNA-CNB
interaction region

PxPxIT 1/6
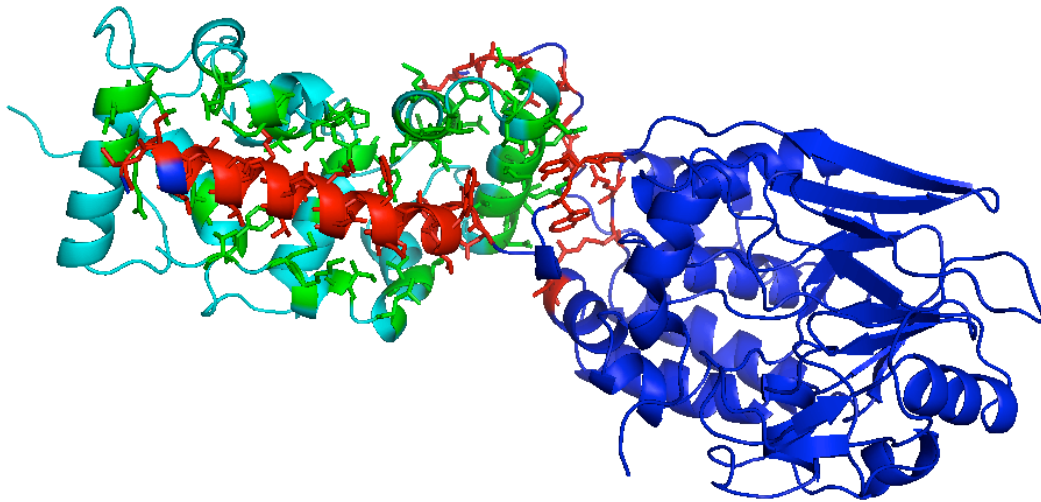
**Figure 7**
InterProSurf

**Figure 8**
PPI-PRED

**Figure 9**
Whiscey

**Figure 10**
Potential new interaction patches

Unverified region #1

| Promate | SPPIDER | mPPISP |
|---------|---------|--------|
| P204 | T249 | **Y224** |
| E205 | **H250** | G255 |
| **Y224** | T252 | Y258 |
| D229* | V253 | M483 |
| **H250** | G255 | P484 |
| **C256** | **C256** | |

Unverified
Region #2
cPPISP

L190
L209
N272
N273

Unverified
region #3
cPPISP

N77
L215
D216

References:

Bradford J. and Westhead R. (2005). Improved prediction of potein-protein binding sitesusing a support vector machines approach. Bioinformatics 21(8), 1487-1494.

Chen H. and Zhou H. (2005). Prediction of Interface Residues in Protein-Protein Complexes by a Consensu Neural Network Method: Test Against NMR Data. Proteins: Structure, Function and Boinformatics 61, 21-35

Headd J., Ban A., Brown P., Edelsbrunner H., Vaidya M., Rudolph J. (2007). Protein-Protein Interfaces: Properties, Preferences, and Projections. Journal of Proteome research 6, 2576-2586.

Liang S., Zhang C., Liu S. and Zhou Y. (2006). Protein inding site prediction using n empirical scoring function. Nucleic Acids Research 34(13), 3698-3707.

Mandell D. and Kortemme T. (2009). Computer-aided design of functional protein interactions. Nature chemical biology 5(11), 797-807.

Negi S., Schein C., Oezguen N., Power T., and Braun W. (2007). InterProSurf: a web server for predicting interacting sites on protein surfaces. Bioinformatics 23(24), 3397-3399.

Neuvirth H., Raz R., and Schreiber G. (2004). ProMate: A Structure Based Prediction Program to Identify the Location of Protein-Protein Binding Sites. J. Mol. Iol. 338, 181-199.

Porollo A., and Meller J. (2007). Prediction-Based Fingerprints of Protein-Protein Interactions. Proteins: Sturcture, Function, and Bioinformatics 66, 630-645,

Qin S. and Zhou H. (2007). A holistic approach to protein docking. Proteins 69,743-749.

Qin S.  and Zhou H. (2007b). Meta-PPIS: a meta web server for protein-protein interaction site prediction. Bioinformatics 23(24), 3386-3387.

Via A., Ferre F., Brannetti B., and Helmer-Citterich M. (2000). Protein surface similarities: a survey of methods to describe and compare protein surfaces. Cell. Mol. Life Sci. 57, 1970-1977.

De Vries S., van Dijk A. and Bonvin A. (2006). WHISCY: What Information Does Surface Conservation Yield? Applicatoin to Data-Driven Docking. Proteins: Sturcture, Function, and Bioinformatics 63, 479-489.

Yan C, Wu F, Jernigan RL, Dobbs D, Honavar V. (2008). Characterization of protein-protein interfaces. Protein J. Jan;27(1):59-70.;

Yan C, Dobbs D, Honavar V. A two-stage classifier for identification of protein-protein interface residues. Bioinformatics. 2004 Aug 4;20 Suppl 1:i371-8.).

Zhou H. and Qin S. (2007) Interaction-site prediction for protein complexes: a critical assessment. Bioniformatics 23(17), 2203-2209.