# GalaxyHomomer: a web server for protein homo-oligomer structure prediction from a monomer sequence or structure

**Minkyung Baek[1], Taeyong Park[1], Lim Heo[1], Chiwook Park[2] and Chaok Seok[1,*]**

[1]Department of Chemistry, Seoul National University, Seoul 151-747, Korea and [2]Department of Medicinal Chemistry and Molecular Pharmacology, Purdue University, West Lafayette, IN 47907, USA

## ABSTRACT

**Homo-oligomerization of proteins is abundant in nature, and is often intimately related with the physiological functions of proteins, such as in metabolism, signal transduction or immunity. Information on the homo-oligomer structure is therefore important to obtain a molecular-level understanding of protein functions and their regulation. Currently available web servers predict protein homo-oligomer structures either by template-based modeling using homo-oligomer templates selected from the protein structure database or by *ab initio* docking of monomer structures resolved by experiment or predicted by computation. The GalaxyHomomer server, freely accessible at http://galaxy.seoklab.org/homomer, carries out template-based modeling, *ab initio* docking or both depending on the availability of proper oligomer templates. It also incorporates recently developed model refinement methods that can consistently improve model quality. Moreover, the server provides additional options that can be chosen by the user depending on the availability of information on the monomer structure, oligomeric state and locations of unreliable/flexible loops or termini. The performance of the server was better than or comparable to that of other available methods when tested on benchmark sets and in a recent CASP performed in a blind fashion.**

## INTRODUCTION

A large fraction of cellular proteins self-assemble to form symmetric homo-oligomers with distinct biochemical and biophysical properties (1–3). For example, ligand-binding sites or catalytic sites are located at oligomer interfaces in many proteins (4–6), and oligomerization is often necessary for effective signal transduction through membrane receptor proteins (7,8) and selective gating of channel proteins (9). Therefore, knowledge of the homo-oligomer structure is essential for understanding the physiological functions of proteins at the molecular level and for designing molecules that regulate the functions.

Methods for predicting the protein homo-oligomer structure can be divided into two categories: those that use templates selected from the protein structure database and others that dock monomer structures *ab initio*, without using template information. Usually, template-based methods require a sequence as input, whereas docking methods require a monomer structure as input. The latter requirement can be more restrictive for the user if the monomer structure has to be predicted by another method, but it may be preferred if an experimentally resolved monomer structure is available. It is generally expected that template-based methods produce more accurate predictions under a situation in which similar proteins forming oligomers exist in the structure database. Docking methods may be more useful when proper oligomer templates are not available but the monomer structure is reliable. Several protein–protein docking methods have been reported to date (10–18), and some of these are available as public web servers for predicting homo-oligomer structures. M-ZDOCK (13) and GRAMM-X (15), which use *ab initio* docking based on fast Fourier transformation (FFT), are two such examples. The oligomeric state must be provided as input in these servers. However, relatively few web servers that use template-based methods have been reported. ROBETTA (19,20) and SWISS-MODEL (21) are two web servers that predict the homo-oligomer structure from an amino acid sequence. GalaxyGemini (22) predicts the homo-oligomer structure from a monomer structure. These servers predict the oligomeric state automatically. Depending on the availability of information on the oligomeric state, the user may or may not prefer to specify the oligomeric state. Here, we introduce a new web server called GalaxyHomomer that predicts the homo-oligomer structure from either the

*To whom correspondence should be addressed. Tel: +82 2 880 9197; Fax: +82 2 889 1568; Email: chaok@snu.ac.kr
Present address: Lim Heo, Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI 48824, USA.

amino acid sequence or from the monomer structure. The oligomeric state may or may not be specified by user. The server can perform both template-based oligomer modeling and *ab initio* docking. It returns five model structures and automatically decides how many models are generated by which method depending on the existence of proper oligomer templates.

Oligomer structures predicted by template-based methods may have errors due to sequence differences between the target and template proteins. Those predicted by docking methods may have inaccuracy if structural change of the monomer induced by oligomerization is not considered. In the previous CASP experiment conducted in 2014 in collaboration with CAPRI, we showed that such errors in predicted oligomer structures could be reduced by re-modeling inaccurately predicted loops or termini and by relaxing the overall structure (23). GalaxyHomomer incorporates such state-of-the-art model refinement methods to improve the accuracy of homo-oligomer models generated by both template-based modeling and *ab initio* docking.

According to the assessment of the recent blind prediction experiment CASP12 conducted in 2016, GalaxyHomomer, participated as 'Seok-assembly', ranked second among the servers participated in the assembly category. When we tested GalaxyHomomer on 136 targets from PISA benchmark set, 47 targets from a membrane protein set, 20 targets from CASP11 experiments and 89 targets from CAMEO protein structure prediction category, it showed a performance better than or comparable to that of other available homo-oligomer structure prediction methods.

## THE GALAXYHOMOMER METHOD

### Overall procedure

The overall pipeline of GalaxyHomomer is presented in Figure 1. Either a sequence or structure (experimental or predicted structure) of the monomer can be provided as input. If the oligomeric state is not specified by the user, possible oligomeric states are predicted first. Five homo-oligomer structures with the given oligomeric states are then generated by template-based modeling and *ab initio* docking. Oligomer templates required by template-based modeling are detected based only on the sequence as well as with additional structure information. The models are further refined by loop/terminus modeling using GalaxyLoop (24–26) and by overall relaxation using GalaxyRefineComplex (27).

### Prediction of the oligomeric state

Possible oligomeric states are predicted from the input sequence by a similarity-based method as follows. First, HHsearch (28) is run in the local alignment mode to detect proteins that are similar to the target in the protein structure database 'pdb70', with a maximum mutual sequence identity of 70%. The oligomeric states of the database proteins were assigned according to the biological units described in 'REMARK 350'. Second, the proteins are re-ranked by a score $S$, which combines the HHsearch sequence score and HHsearch secondary structure score (29). Next, the $S$ scores of the proteins in the same oligomeric states are summed for
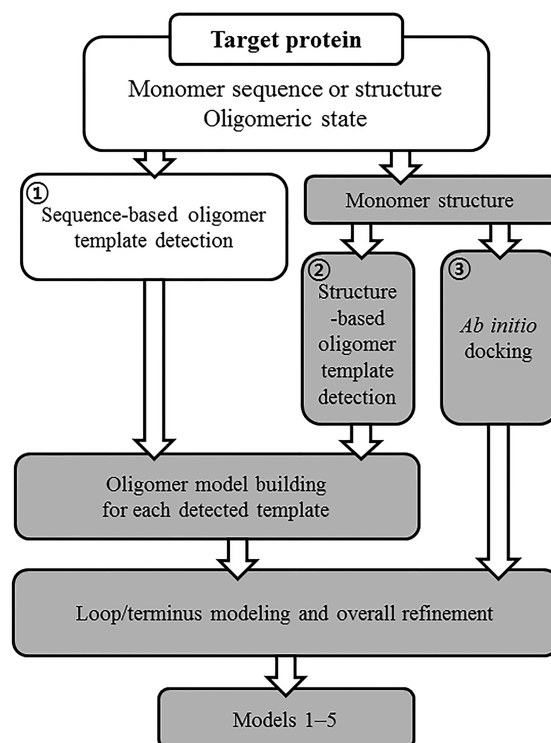


**Figure 1.** Flowchart of the GalaxyHomomer algorithm. The homo-oligomer structure prediction methods based on sequence similarity, structure similarity and *ab initio* docking are attempted in the order in which they are numbered until five homo-oligomer models are generated. When the monomer structure is given as input, only shaded procedures are executed.

the top 100 proteins, and the ratios of different oligomeric states are determined in proportion to the $S$ sums. Finally, oligomeric states for five models are assigned according to the oligomeric state ratios.

### Template-based oligomer modeling

The same top 100 proteins described above are considered as candidates for oligomer templates. If a sequence is provided as input, up to five proteins are selected as templates based on the ranking of $S$ among those with $S > 0.2$ times the highest $S$ overall and those $> 0.7$ times the highest $S$ for the given oligomeric state. If the number of detected templates using this sequence-based method is less than five, additional templates are selected using the monomer structure predicted by the template-based modeling program GalaxyTBM (29). Structure-based templates are selected according to the ranking of $S$ among those with monomer structures similar to the given monomer structure (TM-score calculated using TM-align (30) $> 0.5$) and in the given oligomeric state. If a structure is provided as input, only the structure-based template detection is used with the monomer structure provided by the user.

For each oligomer template detected by the sequence-based method, an oligomer structure is built using the in-house model-building program GalaxyCassiopeia, a component of the most recent version of GalaxyTBM (29).

**Table 1.** Performance comparison of homo-oligomer structure prediction methods in terms of the CAPRI accuracy criterion

| Benchmark set | Prediction methods | Input | Up to 5 models[a] | Top 1 model[a] |
|---|---|---|---|---|
| PISA (136 targets)[b] | GalaxyHomomer | Sequence | 62/5***/38** | 57/3***/39** |
|  | HH+MODELLER[c] | Sequence | 61/3***/38** | 45/1***/26** |
| Membrane proteins (47 targets)[b] | GalaxyHomomer | Sequence | 19/1***/14** | 19/1***/9** |
|  | HH+MODELLER[c] | Sequence | 18/0***/6** | 14/0***/4** |
| CASP11 (20 targets)[b] | GalaxyHomomer | Sequence | 12/0***/8** | 12/0***/5** |
|  | HADDOCK | Structure | 14/0***/10** | 13/0***/9** |
|  | ClusPro | Structure | 14/0***/7** | 10/0***/5** |
|  | BAKER-ROSETTASERVER | Sequence | 9/0***/8** | 9/0***/7** |
|  | SwarmDock | Structure | 9/0***/3** | 8/0***/3** |
|  | GalaxyGemini[d] | Structure | Not available | 7/0***/5** |
|  | GRAMM-X | Structure | 5/0***/1** | 3/0***/1** |
| CAMEO (89 targets) | GalaxyHomomer | Sequence | 44/6***/25** | 35/3***/25** |
|  | Robetta | Sequence | 28/4***/17** | 26/4***/15** |
|  | SWISS-MODEL[d] | Sequence | Not available | 23/3***/16** |

[a]Data represent the numbers of targets for which the best of up to five predicted models were of acceptable or higher/high accuracy (***) and medium accuracy (**); values for model 1 are shown.
[b]Oligomeric state of target protein is given as an input.
[c]Up to five homo-oligomer models were generated by MODELLER based on the templates detected by HH-search.
[d]Data for up to five models were not provided for GalaxyGemini and SWISS-MODEL because they generated only single models.

GalaxyCassiopeia builds models from the sequence alignment and template structure by the VTFM optimization used in MODELLER (31) but with FACTS solvation free energy (32), knowledge-based hydrogen bond energy (33) and dipolar-DFIRE (34) in addition to molecular mechanics bonded and non-bonded energy terms and template-derived restraints. For each template detected by the structure-based method, an oligomer structure is built by superimposing the monomer structure onto the oligomer template.

### *Ab initio* docking

If less than five oligomer templates are detected by the two template detection methods described above, the remaining homo-oligomer models with the given oligomeric states are generated using the in-house *ab initio* docking program GalaxyOligoTongDock. This docking program predicts homo-oligomer structures from the monomer structure using the grid-based FFT docking method of M-ZDOCK (13) implemented in-house. Only $C_n$-symmetry is considered, and D-symmetry is not supported. The top 200 homo-oligomer structures generated by FFT are clustered using NMRCLUST (35), and the clusters are ranked according to the cluster size. From each of the highest ranking clusters, the highest-score structure is selected.
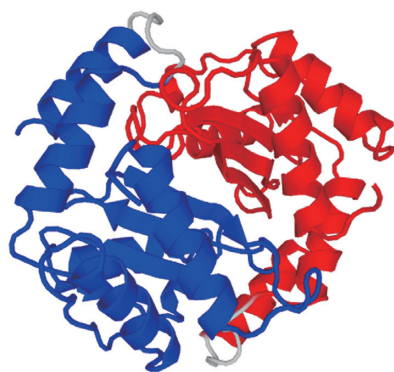
### Structure refinement

Less reliable loop or terminal regions are re-modeled using GalaxyLoop (24–26) considering symmetry of the homo-oligomer structure for the first model for those regions predicted to be unreliable if a sequence is provided as input, and for all five models for user-specified regions if a structure is provided as input. GalaxyRefineComplex (27) is subsequently run to further relax the overall structure. The user can run additional refinement jobs by clicking the 'Submit' button in the results table on the output page.

### Performance of the method

The GalaxyHomomer server was tested on 25 targets in CASP12 in a blind fashion, and this server, named 'Seok-assembly', ranked second among the servers participated in the assembly category (http://www.predictioncenter.org/casp12/). In CASPs, the oligomeric state is provided by the organizers. The server was also tested on three benchmark sets for which the oligomeric state is given as input (136 homo-oligomer proteins from the PISA benchmark set (36), 47 homo-oligomer membrane proteins compiled from the PDB (Supplementary Data) and 20 homo-oligomer proteins among the targets of CASP11 held in 2014 in collaboration with CAPRI (18)) and on a set for which the oligomeric state is not provided as input (89 homo-oligomer proteins among CAMEO (37) targets released from 13 August 2016 to 11 November 2016). In these tests, the performance of GalaxyHomomer was better than or comparable to that of other methods for which performance data are available for the sets in terms of the CAPRI accuracy criterion, as summarized in Table 1. Note that some methods take only the structure as input. The CAPRI criterion reflects the biological relevance of the model structures, and model qualities are classified as high (***), medium (**), acceptable (*) and incorrect considering the ligand root mean-square deviation (L-RMSD) and interface RMSD (I-RMSD) from the experimental structure and the fraction of predicted native contacts ($F_{nat}$) (38). See Supplementary Data for details on the benchmark tests.

It has to be noted that GalaxyHomomer does not consider the lipid bilayer environment of membrane proteins explicitly in terms of energy or geometry during energy-based optimization and docking. However, the results on membrane proteins in Table 1 are quite promising, implying that membrane environment was effectively taken into account in an implicit manner by using the database structures of membrane proteins as templates. GalaxyHomomer showed better performance than GalaxyGemini (22), a previous homo-oligomer structure prediction server developed by us, on the CASP11 benchmark set, as summarized in Ta-

**View in PV** [Model 1] [Model 2] [Model 3] [Model 4] [Model 5]

**Download** [Model 1] [Model 2] [Model 3] [Model 4] [Model 5]

**Template-based Oligomer Modeling Results (Sequence-based)**

| Model No | Oligomer template | Number of subunits | Interface area | Sequence identity | Loops/termini that may be refined | Re-submit for refinement |
|---|---|---|---|---|---|---|
| 1 | 3brc | 2-mer | 2122.0 | 44.6 | 15-21 | Submit |

**Template-based Oligomer Modeling Results (Structure-based)**

| Model No | Oligomer template | Number of subunits | Interface area | Sequence identity | Structure similarity (TM-score) | Loops/termini that may be refined | Re-submit for refinement |
|---|---|---|---|---|---|---|---|
| 2 | 1zwy | 2-mer | 798.0 | 16.6 | 0.5682 | 59-74 | Submit |
| 3 | 2car | 2-mer | 209.9 | 10.2 | 0.5230 | 59-74 | Submit |
| 4 | 3tqu | 2-mer | 2193.4 | 9.6 | 0.5123 | 59-74 | Submit |

**Ab initio Docking Results**

| Model No | Number of subunits | Interface area | Docking score | Loops/termini that may be refined | Re-submit for refinement |
|---|---|---|---|---|---|
| 5 | 2-mer | 628.9 | 625.491 | 59-74 | Submit |

**Download**

- 5 homo-oligomer models [DOWNLOAD]

**Figure 2.** An example output page of GalaxyHomomer. Five generated models are visualized using the JavaScript Protein Viewer. The models can be downloaded in PDB format. Additional information such as the number of subunits, interface area, information on templates and *ab initio* docking score is provided in the tables.

ble 1. The difference in the performance is mainly due to the cases in which predicted monomer structures are not accurate enough. In such cases, oligomer structures built directly from the sequence using sequence-based templates (method 1 in Figure 1) tended to be more accurate than those obtained by superimposing the predicted monomer structures on the structure-based templates (method 2 in Figure 1) GalaxyGemini builds oligomer models using only method 2. Additional model refinement performed by GalaxyHomomer also improved the model accuracy.

## THE GALAXYHOMOMER SERVER

### Hardware and software

The GalaxyHomomer server runs on a cluster of 12 Linux servers of 2.33-GHz Intel Xeon 8-core processors. The web application uses the Python programming language and the MySQL database. The whole GalaxyHomomer pipeline is implemented using Python. The model building, *ab initio* docking and refinement methods are implemented as part of the GALAXY program package (39) written in Fortran

90. The JavaScript Protein Viewer (http://biasmv.github.io/pv/) is used for visualization of the predicted models.

### Input and output

The required input is a protein monomer sequence in FASTA format or a protein monomer structure in PDB format. The number of residues in the input file is limited to 1000 for computational efficiency. Users can provide additional information such as the oligomeric state and locations of unreliable/flexible loops or termini of the input structure. Usual run time is 6–12 h, but it depends heavily on the homo-oligomer size and the input type, as shown in Supplementary Figure S2. Five homo-oligomer model structures are visualized and available for download in PDB format. Detailed prediction results, including the number of subunits, interface area, information on templates and *ab initio* docking score, are reported in the tables (Figure 2).

## CONCLUSION

The GalaxyHomomer server predicts the homo-oligomer structure of a target protein from a sequence or monomer structure. It performs both template-based modeling and *ab initio* docking, and adopts additional model refinement that can consistently improve model quality. The server provides different options that can be chosen by the user depending on the availability of information on monomer structure, oligomeric state and locations of unreliable/flexible loops or termini. By combining additional refinement based on loop modeling and overall structure refinement, GalaxyHomomer may generate more precise homo-oligomer models that can be useful for further applications such as for drug design targeting protein homo-oligomer interfaces.

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## FUNDING

## REFERENCES

1. Andre,I., Strauss,C.E., Kaplan,D.B., Bradley,P. and Baker,D. (2008) Emergence of symmetry in homooligomeric biological assemblies. *Proc. Natl. Acad. Sci. U.S.A.*, **105**, 16148–16152.
2. Goodsell,D.S. and Olson,A.J. (2000) Structural symmetry and protein function. *Ann. Rev. Biophys. Biomol. Struct.*, **29**, 105–153.
3. Poupon,A. and Janin,J. (2010) Analysis and prediction of protein quaternary structure. *Methods Mol. Biol.*, **609**, 349–364.
4. Snijder,H.J., Ubarretxena-Belandia,I., Blaauw,M., Kalk,K.H., Verheij,H.M., Egmond,M.R., Dekker,N. and Dijkstra,B.W. (1999) Structural evidence for dimerization-regulated activation of an integral membrane phospholipase. *Nature*, **401**, 717–721.

5. Ali,A., Bandaranayake,R.M., Cai,Y., King,N.M., Kolli,M., Mittal,S., Murzycki,J.F., Nalam,M.N., Nalivaika,E.A., Ozen,A. *et al.* (2010) Molecular basis for drug resistance in HIV-1 protease. *Viruses*, **2**, 2509–2535.

6. Pidugu,L.S.M., Mbimba,J.C.E., Ahmad,M., Pozharski,E., Sausville,E.A., Emadi,A. and Toth,E.A. (2016) A direct interaction between NQO1 and a chemotherapeutic dimeric naphthoquinone. *BMC Struct. Biol.*, **16**, 1.

7. Heldin,C.H. (1995) Dimerization of cell surface receptors in signal transduction. *Cell*, **80**, 213–223.

8. Stock,J. (1996) Receptor signaling: dimerization and beyond. *Curr. Biol.*, **6**, 825–827.

9. Clarke,O.B. and Gulbis,J.M. (2012) Oligomerization at the membrane: potassium channel structure and function. *Adv. Exp. Med. Biol.*, **747**, 122–136.

10. Dominguez,C., Boelens,R. and Bonvin,A.M. (2003) HADDOCK: a protein-protein docking approach based on biochemical or biophysical information. *J. Am. Chem. Soc.*, **125**, 1731–1737.

11. Gray,J.J., Moughon,S., Wang,C., Schueler-Furman,O., Kuhlman,B., Rohl,C.A. and Baker,D. (2003) Protein-protein docking with simultaneous optimization of rigid-body displacement and side-chain conformations. *J. Mol. Biol.*, **331**, 281–299.

12. Comeau,S.R., Gatchell,D.W., Vajda,S. and Camacho,C.J. (2004) ClusPro: an automated docking and discrimination method for the prediction of protein complexes. *Bioinformatics*, **20**, 45–50.

13. Pierce,B., Tong,W. and Weng,Z. (2005) M-ZDOCK: a grid-based approach for Cn symmetric multimer docking. *Bioinformatics*, **21**, 1472–1478.

14. Schneidman-Duhovny,D., Inbar,Y., Nussinov,R. and Wolfson,H.J. (2005) PatchDock and SymmDock: servers for rigid and symmetric docking. *Nucleic Acids Res.*, **33**, W363–W367.

15. Tovchigrechko,A. and Vakser,I.A. (2006) GRAMM-X public web server for protein-protein docking. *Nucleic Acids Res.*, **34**, W310–W314.

16. Macindoe,G., Mavridis,L., Venkatraman,V., Devignes,M.D. and Ritchie,D.W. (2010) HexServer: an FFT-based protein docking server powered by graphics processors. *Nucleic Acids Res.*, **38**, W445–W449.

17. Torchala,M., Moal,I.H., Chaleil,R.A., Fernandez-Recio,J. and Bates,P.A. (2013) SwarmDock: a server for flexible protein-protein docking. *Bioinformatics*, **29**, 807–809.

18. Lensink,M.F., Velankar,S., Kryshtafovych,A., Huang,S.Y., Schneidman-Duhovny,D., Sali,A., Segura,J., Fernandez-Fuentes,N., Viswanath,S., Elber,R. *et al.* (2016) Prediction of homoprotein and heteroprotein complexes by protein docking and template-based modeling: a CASP-CAPRI experiment. *Proteins*, **84**(Suppl. 1), 323–348.

19. Kim,D.E., Chivian,D. and Baker,D. (2004) Protein structure prediction and analysis using the Robetta server. *Nucleic Acids Res.*, **32**, W526–W531.

20. DiMaio,F., Leaver-Fay,A., Bradley,P., Baker,D. and Andre,I. (2011) Modeling symmetric macromolecular structures in Rosetta3. *PLoS One*, **6**, e20450.

21. Biasini,M., Bienert,S., Waterhouse,A., Arnold,K., Studer,G., Schmidt,T., Kiefer,F., Gallo Cassarino,T., Bertoni,M., Bordoli,L. *et al.* (2014) SWISS-MODEL: modelling protein tertiary and quaternary structure using evolutionary information. *Nucleic Acids Res.*, **42**, W252–W258.

22. Lee,H., Park,H., Ko,J. and Seok,C. (2013) GalaxyGemini: a web server for protein homo-oligomer structure prediction based on similarity. *Bioinformatics*, **29**, 1078–1080.

23. Lee,H., Baek,M., Lee,G.R., Park,S. and Seok,C. (2016) Template-based modeling and ab initio refinement of protein oligomer structures using GALAXY in CAPRI round 30. *Proteins*, **85**, 399–407.

24. Lee,J., Lee,D., Park,H., Coutsias,E.A. and Seok,C. (2010) Protein loop modeling by using fragment assembly and analytical loop closure. *Proteins*, **78**, 3428–3436.

25. Park,H. and Seok,C. (2012) Refinement of unreliable local regions in template-based protein models. *Proteins*, **80**, 1974–1986.

26. Park,H., Lee,G.R., Heo,L. and Seok,C. (2014) Protein loop modeling using a new hybrid energy function and its application to modeling in inaccurate structural environments. *PLoS One*, **9**, e113811.

27. Heo,L., Lee,H. and Seok,C. (2016) GalaxyRefineComplex: refinement of protein-protein complex model structures driven by interface repacking. *Sci. Rep.*, **6**, 32153.

28. Soding,J. (2005) Protein homology detection by HMM-HMM comparison. *Bioinformatics*, **21**, 951–960.

29. Ko,J., Park,H. and Seok,C. (2012) GalaxyTBM: template-based modeling by building a reliable core and refining unreliable local regions. *BMC Bioinformatics*, **13**, 198.

30. Zhang,Y. and Skolnick,J. (2005) TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids Res.*, **33**, 2302–2309.

31. Sali,A. and Blundell,T.L. (1993) Comparative protein modelling by satisfaction of spatial restraints. *J. Mol. Biol.*, **234**, 779–815.

32. Haberthur,U. and Caflisch,A. (2008) FACTS: fast analytical continuum treatment of solvation. *J. Comput. Chem.*, **29**, 701–715.

33. Kortemme,T., Morozov,A.V. and Baker,D. (2003) An orientation-dependent hydrogen bonding potential improves prediction of specificity and structure for proteins and protein-protein complexes. *J. Mol. Biol.*, **326**, 1239–1259.

34. Yang,Y. and Zhou,Y. (2008) Specific interactions for ab initio folding of protein terminal regions with secondary structures. *Proteins*, **72**, 793–803.

35. Kelley,L.A., Gardner,S.P. and Sutcliffe,M.J. (1996) An automated approach for clustering an ensemble of NMR-derived protein structures into conformationally related subfamilies. *Protein Eng.*, **9**, 1063–1065.

36. Ponstingl,H., Kabir,T. and Thornton,J.M. (2003) Automatic inference of protein quaternary structure from crystals. *J. Appl. Crystallogr.*, **36**, 1116–1122.

37. Haas,J., Roth,S., Arnold,K., Kiefer,F., Schmidt,T., Bordoli,L. and Schwede,T. (2013) The protein model portal–a comprehensive resource for protein structure and model information. *Database*, **2013**, bat031.

38. Lensink,M.F. and Wodak,S.J. (2010) Docking and scoring protein interactions: CAPRI 2009. *Proteins*, **78**, 3073–3084.

39. Ko,J., Park,H., Heo,L. and Seok,C. (2012) GalaxyWEB server for protein structure prediction and refinement. *Nucleic Acids Res.*, **40**, W294–W297.