

# Predicting the conformational preferences of proteins using a physics-based free energy method

Arijit Roy, Alberto Perez, Ken A. Dill, and Justin L. MacCallum

Laufer Center for Physical and Quantitative Biology  
and Departments of Physics and Chemistry  
Stony Brook University  
Stony Brook, NY 11794-5252.

April 29, 2013

## Abstract

Calculation of free energy differences is of central importance in the simulation of biochemical systems. It is particularly difficult to calculate between pairs of macromolecular conformations as well as a computationally expensive task with existing methods. In this work, confinement approach is used to calculate absolute free energies of biomolecular systems. This method provides two main advantages: it does not require a reaction coordinate or transition path and it is fast to compute. Free energy calculated can be decomposed into a per residue contribution in an approximate way. Per residue free energy allows us to identify the reason behind conformational preferences in biomolecules. Through out the article we show its use in different challenging modeling problems. In particular, we show its use in predicting the conformational preference of chameleon sequences (sequences with high sequence identity and different folds). This sequence dependent conformational preferences and per residue free energy decomposition set the stage for the use of this method in protein design.

[Guys—I recommend we get a better name than confinement. It's totally uninformative. In my view, better would be 'confine-and-configure' or something. Please give it some thought.]

[Ken—I agree that the name isn't great, but I don't think we should change it because: it's already in the literature and we didn't come up with the method. I think we should consider capitalizing it as Confinement Method throughout.]

## 1 Introduction

In some problems of protein science, you want to know the relative stability of a protein's conformation  $A$  compared to its conformation  $B$ . We call this the *difference free energy*. For example, in all esoteric mechanisms, a protein often adopts one conformation when a ligand is bound and another conformation when no ligand is bound<sup>1</sup>. To understand the mechanism requires knowledge of the relative free energies of the protein conformations. Another example is in protein structure prediction, for example as practiced in the Critical Assessment of Structure Prediction (CASP) competition<sup>2</sup>. If you are using some computational model and have predicted two putative native structures,  $A$  and  $B$ , you want to compute which has the lower free energy, in order to know which is the more native-like. In a third example, you may want to know whether mutating a few amino acids in a protein could cause the protein to switch from one stable conformation to another, because that could have important consequences for biological mechanisms and disease. And fourth, sometimes binding a ligand induces a protein from one conformation to another. A quantitative understanding of such *induced-fit* situations requires knowledge of the difference free energy of the protein in the two states. In all these cases, it would be useful to have a computational method that can be used efficiently with atomically detailed physical forcefield models to compute the free energy differences between different given protein conformations.

[TI is actually the basis of the confinement step. We need to be more explicit about what we mean here.] Methods like thermodynamic integration<sup>11, 12</sup> which are successful in alchemical free energy calculation may not be as useful for calculation of conformational free energy. One widely explored

strategy is to use molecular dynamics simulations along some putative reaction coordinate pathway from conformation  $A$  to  $B$ <sup>1,3–8</sup>. [we need many refs here: cite Elber's milestones, Jhih-Wei Chu, and the many other path-sampling approaches.] The free energy along this reaction coordinate can then be determined using methods such as umbrella sampling<sup>10</sup>[add ref where this is used computing free energy differences between conformations] or targeted MD [refs]. Such approaches have several limitations. First, it is necessary to know an efficient reaction pathway from  $A$  to  $B$ . If conformations  $A$  and  $B$  are quite different, then it can be challenging to find such paths. Second, these methods are computationally slow. To get an accurate estimate of the total free energy difference  $\Delta G = G_B - G_A$  requires accurate determinations of the many small free energy differences  $A \rightarrow 1 \rightarrow 2 \dots \rightarrow B$ , so each step requires substantial amounts of sampling. And third, these methods are prone to large errors, because the pathway error is a sum of many errors among the many steps.

[Are all these methods below ones that don't involve paths? Also, are there other key methods we should cite? We need a good scholarly list of all the relevant works.]

The calculation of protein conformational free energy has been successfully attempted by a number of groups<sup>24–32</sup>. Some of these methods like the reference system method<sup>24</sup>, deactivated morphing<sup>25</sup>, the confinement method<sup>28–30</sup> take an alternative strategy for computing difference free energies that does not require knowing a pathway from  $A$  to  $B$ . [please add the several new Karplus papers from 2012–2013.].

Here, we adopt the confinement method of Tyka et al<sup>28</sup> and Cecchini et al.<sup>29</sup>, which is based on the thermodynamic cycle shown in Figure 4. We start with an ensemble of related conformations,  $A$ . We then restrict  $A$  to a much “tighter” ensemble  $A^*$  by applying position restraints in a series of MD simulations. Similarly, a second ensemble  $B$  is restricted to  $B^*$ . We then compute the free energy between  $A^*$  and  $B^*$  using either normal mode analysis or the quasiharmonic method [ref]. This step takes into account both the remaining conformational entropy in each ensemble and the remaining enthalpy difference. In this way, we can calculate the free energy difference between the two end states without needing to define a physical path or reaction coordinate connecting the two states. The confinement method shares some similarities with the ‘confine-and-release’

method for computing ligand binding affinities<sup>13–15</sup> [cite mobley's jctc 07 (ref 216 on our labsite), and mobley's recent j chem phys review, and mobley (ref 199) j chem phys, and others that are relevant.].

[Figure 1 about here.]

Previously, the confinement method has been validated on small model peptides<sup>28,29</sup>. Here, we show that the confinement method can also be applied successfully to larger systems, such as the series of chameleon proteins designed by Orban and co-workers<sup>34–36;38</sup>, which can switch between two completely different folds with only slight changes in sequence. We also show that our computed difference free energies are useful in evaluating the quality of CASP target protein predictions. This could be of significant value for ultimately improving the energetics in protein-structure models, not just the structures. And, finally, we show that we can approximately decompose difference free energies into individual amino acid components. This offers the opportunity for diagnosing the structural basis for difference free energies, which can be useful for interpreting biological mechanisms.

## 2 Results and Discussion

### 2.1 The confinement method succeeds at some basic consistency checks

First, we validated that our implementation of the confinement method produces results compatible with previous calculations reported in the literature. The method has previously been applied to a 16 amino acid residue  $\beta$ -hairpin from protein G, known as BHP<sup>29</sup>. We calculated the free energy difference between the native conformation—called bhp1, which has a two-stranded  $\beta$ -sheet—and a non-native conformation—called bhp3, which has a three-stranded  $\beta$ -sheet. Our confinement calculation shows that bhp1 is more stable by 1.7 kcal/mol, consistent with 4  $\mu$ s equilibrium simulations showing that bhp1 is more favorable configuration by 1.8 kcal/mol, and in agreement with previous calculations<sup>29</sup>.

Second, we looked at 6 target proteins from the CASP9 experiment [ref]. For each target, we examined up to 5 submitted models. We computed

the difference free energy between the true native and the best model. Figure 4 and Table 1 show that, in 5 out of 6 cases, the confinement method assigns a lower free energy to the experimentally determined structure than to any of the decoys. Other well-known discriminators can also successfully tell native structures from computer generated models<sup>40</sup> [add rosetta holes, rosetta ff, dfire]; it is just an independent useful validation that the confinement calculations make sense.

[Arijit, please add GDT of best decoy to this table.]

[Table 1 about here.]

## 2.2 The confinement method correctly predicts the structures of chameleon sequences

We also tested the confinement method predictions for difference free energies on the chameleon sequences of Orban et al.<sup>34-39</sup>. These are instances in which two highly similar sequences fold into remarkably different structures. Orban and co-workers have designed a protein-G-like sequence of 56-residues that is marginally stable in one of two possible folds. By mutating key residues in this sequence they are able to stabilize one fold or the other (see Figure 4). We refer to the  $4\beta + \alpha$  structure as the  $\beta$  conformation, and we refer to the  $3\alpha$  structure as the  $\alpha$  conformation. One pair of sequences (GA88/GB88) is 88 percent identical in sequence, differing in seven positions. Another pair (GA95/GB95) is 95 percent identical, differing in three positions. Accurately predicting the structural preferences of these structures presents a serious challenge for computational methods<sup>34-39</sup>. Do any of these citations try to do any modeling or comment about how challenging this would be?

[Figure 2 about here.]

We initially approached this problem by making a model of each sequence with the same backbone structure as its partner chameleon sequence. For example, we took the sequence of GA88 and built a model with the same overall structure as GB88. We then used the confinement method to assess the free energy difference between the experimentally determined structure

of GA88 and the model (with the GA88 sequence and the GB88 structure). The confinement method was able to predict the conformational preferences correctly for all four sequences (data not shown). This is, however, not surprising as it is often easy to distinguish computational models from native structures<sup>19,40</sup>. To avoid this potential problem, we instead computed relative free energy of two different computer generated models for each sequence. One model is based on the  $\alpha$  structure and the other on the  $\beta$  structure (see Supporting Information for details on the modeling procedure). This is a much more realistic test of the confinement method's ability to accurately calculate the difference free energy.

[ It's not clear what this test is and why it's better, and why there are only 5 sequences. Please clarify.] [There are five pairs of sequences that can be arranged into 3 different pairings. Arijit has removed the 5th sequence.]

Now there are 4 sequences as they have the experimentally available structure. The fifth sequence is discussed at the end of the paragraph.

[Ken: it's not that hard to tell an experimental structure from a computer generated model. Almost all of the knowledge-based scoring functions can do this a substantial fraction of the time. Comparing a near-native model with another model is a much more challenging problem.]

Figure 4 shows that the confinement method identifies the correct structure for all four sequences. One hypothesis is that such structural transitions require states with diminished stability<sup>37</sup>. It is believed that if the free energy of the native state and the alternative state are within a range of around 5 kcal/mol, then it is possible for the native state to be destabilized relative to the alternative state with only small changes in sequence<sup>35-37</sup>. The calculated free energy differences range from around 3.5 to 5.0 kcal/mol, which is consistent with this hypothesis. In a more recent study<sup>38</sup>, the amino acid residue at position 45 (Tyr for  $\beta$  and Leu for  $\alpha$ ) was found to be important for switching between  $\alpha$  and  $\beta$  conformation. This inspired us to introduce another mutation at this position, Y45A, which we refer to as GA98. Our calculations predict that this mutation shifts the equilibrium to the  $\alpha$  conformation, which is now more stable than the  $\beta$  conformation by 3.8 kcal/mol. Although this result has not yet been confirmed experimentally, it is consistent with the previously observed effect of Y45L.

## 2.3 Per-residue free energy calculations can identify the mechanistic detail behind conformational preferences

To better understand the mechanism behind these chameleon proteins, we decomposed the calculated free energy into per-residue contributions in an approximate way<sup>28</sup>. We decompose each confinement step ( $\Delta G_{A,A^*}$ ,  $\Delta G_{B,B^*}$ ) into per-residue contributions. We can also decompose the remaining enthalpy in the confined ensemble into per-residue contributions. However, the method is approximate because we do not include the residual conformational entropy from the normal mode or quasiharmonic analysis steps. Computing per-residue contributions helps us to identify important residues that stabilize a particular conformation. The per residue free energy,  $\Delta\Delta G((\beta) - (\alpha))$  is shown in Figure 4.

[Figure 3 about here.]

Although the overall free energy difference between the two structures is small (within 3–5 kcal/mol), individual residues can show marked preferences for being in either  $\alpha$  or  $\beta$  conformation. Such differences can be understood by looking in detail at the local environment for those residues. For example, the region around residue 7 forms a random coil in the  $\alpha$  structure, whereas it is forming a beta sheet structure in the  $\beta$  structure. These residues strongly favour the locally well packed and hydrogen bonded environment found in the  $\beta$  sheet. Overall, favoring  $\alpha$  or  $\beta$  structure is a delicate balance, were the relative global free energy difference is small and the contributions of different per residue tendencies balance out. It is therefore very likely that by changing key residue preferences the global fold preferences can be changed.

Three residues are mutated between GA95 and GB95, at positions 20, 30, and 45. In GA95, these residues (L20, I30 and L45) stabilize the  $\alpha$  structure (compare the upper and lower panels on Figure 4). In GB95, two of these residues (F30 and Y45) favor the  $\beta$  structure, because they have large solvent exposed surface areas in the  $\alpha$  structure, but are more buried in the  $\beta$  structure. Additionally Y45 forms a hydrogen bond with D47 in the  $\beta$  structure. On the other hand, residue A20 from GB95 still favors the  $\alpha$  structure, although less strongly than L20 in GA95.

The  $\alpha$  and  $\beta$  sequences are nearly identical, so there are some common features observed for all sequences. The experimental observations classified the protein into two parts: Amino acids 9–51 are fully structured in both folds, whereas residues 1–8 and 52–56 are unstructured in the  $\alpha$  fold, but form  $\beta$ -strands in the  $\beta$  fold. Most of the amino acid residues in the region 1–9 have negative per-residue free energies, which means that these residues favor the  $\beta$  structure. Roles of some other important residues in stabilizing either the  $\alpha$  or  $\beta$  conformation are summarized in Figure S1.

In addition to the direct effects of the mutation, there are also indirect effects due to small perturbations in the environment around the mutations. For example, the L20A mutation causes a slight repacking around residue 20. This causes large changes in the per-residue free energies of nearby residues T25 and A26. However, the changes for these two residues have opposite signs and nearly cancel.

These per-residue free energy decompositions provide a great deal of insight into the driving forces behind protein folding and conformational change. We believe that such calculations may also be useful for protein design—designing proteins with specific structures and functions.

## 2.4 The confinement method is a useful tool for structure prediction

We have tested the ability of the confinement method to act as a “meta-predictor” for structure prediction. Here, the task is to correctly identify the most accurate models out of a set of “decoys” generated by different methods during the CASP experiment. CASP is a blind test in which different groups apply methods to predict the 3-dimensional structure of proteins from their sequences. Each group is allowed to submit five possible structures, which they are supposed to rank from best to worst.

We have performed two experiments centered on CASP. In the first experiment, we tried to rank-order predictions for several targets. For each target, the predictions were either produced by a single group—presumably using the same method for each prediction, or were produced by several different groups—using different methods. The goal of this experiment is to determine if the physics-based confinement method can correctly identify more

native-like structures as having lower free energies. Our second experiment was to see if the confinement method can identify structures that are missed by other meta-prediction servers. Most successful meta-prediction servers are based on the idea of consensus: if many different prediction methods produce similar results, then that is probably a correct prediction [ref]. This is often a powerful heuristic, but it can miss cases where there is a very good result that is only predicted by one method. We chose several cases where such structures were missed by the best meta-predictors in CASP and assessed if the confinement method can correctly identify these accurate models.

As is common in the CASP experiment, we assess our results in terms of Global Distance Test Total Score (GDT-TS)<sup>45</sup>, which is a  $C\alpha$  based measure of structural accuracy. It can be understood roughly as the percentage of residues that are correctly positioned in the model (range 0 to 100, higher is better). We did have enough computer model to analyze every model, so the initial models for the test were chosen from a selection of different server groups that have done well in past CASP events.

Overall, the confinement method performs well on our CASP tests. Figure 4 shows that in almost every case, the native structure has the lowest free energy and the best model has the next lowest free energy. The confinement method appears to be a useful tool for ranking structure predictions.

[Figure 4 about here.]

#### 2.4.1 The confinement method can correctly rank-order structure predictions

First, we examined the ability of the confinement method to rank different models for a target that have been generated by a single group using the same methodology for all predictions. We examined two targets: T0559 and T0560 (see Table ?).

[Arijit: Let's move the protein names, organisms, pdb ids, etc to a separate table and just use the CASP ids.]

The first test case is target T0559. The best predictor group for this 69 amino acid target was "BAKER-ROSETTASERVER". To save computer

time, we excluded two models that were very similar to other models that we did include. For this target, the difference free energy computed by the confinement method can be used to accurately rank-order all of the models and the native structure (Figure 4).

Arijit- we could probably make a better visualization of this figure above by putting it onto a toy landscape: x-axis is rmsd, y-axis is DDG, and above each of the 4 points on this graph will be the ribbon diagram and the GDT-TS. I suggest we do that with the fig below too. Let's don't use 4 different colors for the fig below. Let's just show differences in structures in some different color, or show them all as green.]

Ken: I don't mind the 3d plots that are in here now, but I'm not sure I would like something that looks like a free energy landscape.

[Figure 5 about here.]

I'm not sure we need a second figure of this. Move it to the SI?

We performed a similar calculation for target T0560 with two models from the group "Splicer". The remaining three models were discarded as they were too similar to the rest of the models. Again, we can correctly identify the native state and our calculated free energy based ranking matches well with GDT-TS (Figures S?).

[Figure 6 about here.]

Next, we tested the ability of the confinement method to rank models for target T0540 that were produced by different prediction groups. The top models from groups "LTB" (Model 1) and "Mufold" (Model 2) were chosen for analysis. Once again, the free energy based ranking correlates well with the GDT-TS based score (Figure 7).

[Figure 7 about here.]

#### 2.4.2 The per-residue free energy is sensitive to small changes in protein conformation

In section 2.3 we discussed how the per-residue free energy can reveal mechanistic detail behind the conformational preference of a chameleon

sequence. In this section, our aim is to apply the same method and try to understand if the per-residue free energy can help us identify residues which stabilize or destabilize a particular region of a protein. In this case, the two conformations have similar folds with only small changes in localized regions. We chose CASP target T0569 and compared the experimental NMR structure with the best predicted model (GDT-TS=78; predicted by the “Mufold” group). The confinement method predicts that the native structure is more stable by 20 kcal/mol (Figure S2).

The confinement method clearly identifies two hydrophobic residues V59 and I61, which destabilize the predicted model with respect to the crystallographic structure (Figure 8 and Figure S3). The sidechains of these hydrophobic residues are oriented towards the protein hydrophobic core in the native NMR structure but oriented towards the exterior of the protein and are solvent exposed in the model. These residues also form part of a beta-sheet in the experimental structure, but do not form inter-strand hydrogen bonds in the predicted model. There is also a large difference around K76, which forms a salt-bridge with D11 in the predicted model, but not in the experimental structure. This suggests that salt-bridge interactions are too favorable for the combination of force field and implicit solvent model we use, which has been a problem noted in the past [ref carlos].

[Arijit: Figure 8 says crystal, but the text says NMR.]

[Figure 8 about here.]

#### 2.4.3 The confinement method occasionally produces incorrect results

Arijit: let's move Figure 9 to the supporting info and move Figure S4 into the main text.

Despite success for most of the studied systems, there are few failures, specifically when the GDT-TS scores of the compared structures are very close. One example is Target T0538, where we compared the crystal structure with three models (Model 1: “PconsR”—GDT-ts=96; Model 2: “Shell”—GDT-TS=90; Model 3: “FOLDIT”—GDT-TS=86). Contrary to our expectation, the confinement method predicts that Model 1 is more stable than the crystal structure (Figure ?). Per-residue free energy calculations (not shown) show

that despite only small variations at the backbone level, the sidechains are oriented in very different ways (Figure ?), giving rise to large differences in the stabilization of certain residues. In particular, some of the differences arise from different salt bridge patterns and certain flexible polar residues exposed to the surface. This unexpected result shows that the confinement method is very sensitive to local interactions (including sidechain reorientation) and may indicate issues with the forcefield and implicit solvent models used in the calculation.

[Figure 9 about here.]

#### 2.4.4 What can we say about low resolution models?

I think we should consider dropping this whole section. We have a sample size of one and I don't expect that we can generally do well ranking models with a GDT of 44. I think this is a fluke.

So far we have seen that the method is good at predicting preferences when the structures are not very far from the native. But the question remains how far from native can we go and still see that the method produces correct result. In this section we explored this question with models from extracellular domain of the jumping translocation breakpoint protein (pdb id: 2KJX). Most of the group could only generate low resolution models for this CASP9 target (id: T0531). In our comparison, we choose five models by the group MUFOOLD-MD, which was the best performing group for this target with their best model had a GDT\_TS value of 44. The result presented in Figure 10 shows: 1.) native is correctly identified as expected and 2.) Surprisingly there is a high level of correlation between the GDT ranking and the free energy ranking for model 1 and model 3, the rest three structures with GDT\_TS score less than 35 are ordered incorrectly. It is worth to note that models 2 and 3 have the same GDT and very different free energies, meaning that the actual ordering could change a lot<sup>46</sup>. It is encouraging that at least the method can pick out the best model even though it is got a low GDT score: 44.

[Figure 10 about here.]

#### **2.4.5 Can the confinement method perform better quality assessment in protein structure prediction?**

Arijit—is this the case where we chose a structure that was missed by the consensus predictors, if so, we should try to indicate this. The point of this section should be that although consensus prediction strategies do well, they can miss a good answer that is only predicted by a few methods. At least in this case, the confinement method is able to capture this structure that was otherwise missed.

A part of the CASP experiment is dedicated to the quality assessment (QA) of predicted models<sup>44</sup>. Here, predictors were asked to score each model on a scale (known as qmode) from 0 to 1, with higher values corresponding to better models<sup>44</sup>. It will be interesting to know how confinement method perform in quality assessment compare to the other groups in CASP9. We investigated this using couple of CASP Targets. Here, we present a case, where confinement method perform well than the top performing group MUFOLD-WQA<sup>47</sup> of CASP9. We choose two models from CASP target T0538. They were model 3 submitted by PconsR (GDT\_TS = 96, qmode 1 = 0.5434) and model 5 from the MULTICOM-NOVEL (GDT\_TS = 83, qmode 1 = 0.5865). Both are server predicted models and the qmode value presented are from MUFOLD-WQA<sup>47</sup>. We choose these two models as the model by PconsR was the most accurate predicted model for this target, whereas the other model was predicted best by QA test of MUFOLD-WQA. The calculated free energy using confinement method indicate that the model by PconsR is more stable by  $3.9\text{Kcal/Mol}$  which support the GDT\_TS trend. There is no doubt that the consensus approach can predict the model quality in a faster manner. But we expect confinement method can predict the model quality in a relatively expensive but much more accurate way.

### **3 Conclusion**

We have described a computational method called confinement for computing the difference free energy between two conformational ensembles. We show that the difference free energy can be calculated on 100 residue sized proteins, even for large conformational changes. We have demonstrated that it can discriminate the folding preferences of a series of chameleon

proteins. We show that the confinement method can discriminate between the native structure and structure predictions and can identify the best prediction reliably. Perhaps most importantly, we have shown that it can be used to give residue-level insights into what are the dominant structural factors in a protein that are responsible for the difference free energies. The confinement should be useful for protein design, structure prediction, and understanding the mechanism of conformational change.

A key advantage is that this method does not require any reaction coordinate, or sampling a pathway from conformation  $A$  to  $B$ . We have tested the method for structures of proteins having up to 100 amino acids. The computational cost for a 56-residue protein is about 4 hours [on 1 gpu] for 20 ns of confinement [Is 20ns all you need?].

I need 21 (confinement run) x 20 ns for single confinement

[If this is a relevant limitation, let's say something about it.]

I think we need to write the normal mode issue and number of residues that we can handle. If one of you contribute here that will be great.

The conclusion is not the place to talk about how long the method takes or what the technical challenges are. Let's move this to methods section.

## 4 Method

The confinement method has been described in details in ref. by Tyka et al.<sup>28</sup> and Cecchini et al.<sup>29</sup>. The basic approach of the confinement approach is the same in both these papers. A thermodynamic cycle is used to compute the free energy between conformations  $A$  and  $B$ . However, there are some small technical differences between the two approaches. Here we briefly describe the procedure that we used.

1. In the first step, a minimization of  $A$  and  $B$  conformations are performed. These minimized conformations ( $A^{**}$  and  $B^{**}$ ) are the reference conformation of that state.
2. The free energy of confining the ensemble ( $A$  or  $B$ ) to a tighter ensemble ( $A^*$  or  $B^*$ ) is calculated. This is done by gradually confining each

atom to its position in the reference conformation ( $A^{**}$  or  $B^{**}$ ) using a series of progressively stronger position restraints. This is done by running 21 molecular dynamics simulation (each 20 ns long) for each leg of the thermodynamic cycle, where the harmonic restraint force constant was scaled from 0.00005 Kcal/Mol (mostly free) to 81.92 kcal/mol (tightly restrained). The free energy for this step is estimated using a thermodynamic integration approach developed by Tyka et. al.<sup>28</sup>. The confinement free energies calculated this step are denoted as  $\Delta G_{A,A^*}$  and  $\Delta G_{B,B^*}$  in Figure 4.

3. The thermodynamic cycle is closed by calculating the free energy between the final restrained states  $A^*$  and  $B^*$  using normal mode analysis or quasiharmonic analysis. The free energy calculated in this way is shown as  $\Delta G_{A^*,B^*}$  in Figure 4.
4. The full free energy,  $\Delta G_{A,B}$  between the two state A and B is calculated as  $\Delta G_{A,B} = \Delta G_A, B = \Delta G_A, A^* - \Delta G_B, B^* + \Delta G_{A^*,B^*}$ .

All calculations were performed with the Amber 11 suite of programs<sup>41,42</sup> in combination with ff99SB forcefield<sup>49</sup> and generalized born implicit solvent<sup>50</sup>. Interestingly, we extend the method for calculation of per residue free energy in an approximate way. For this purpose, the confinement energy,  $\Delta G_{A,A^*}$  and  $\Delta G_{B,B^*}$  of each residue is calculated in the usual numerical way as described in ref. by Tyka et al.<sup>28</sup>. The internal energy of each residue is calculated using the decomp module of amber from the final restrained trajectory. We call this method approximate as we ignore the entropic contribution from the normal mode or quasiharmonic analysis. However, this contribution to the total free energy is much smaller, which allows us to study the mechanistic details of conformational preference of each residue.

## References

- [1] Elber, R. A Milestoning Study of the Kinetics of an Allosteric Transition: Atomically Detailed Simulations of Deoxy Scapharca Hemoglobin. Biophysical J., 2007, 92, 85-87.

- [2] Moult, J.; Fidelis, K.; Kryshtafovych, A.; Tramontano, A. Critical assessment of methods of protein structure prediction (CASP)-round IX. *Proteins*, 2011, 79, 1-5.
- [3] West, A.M.; Elber, R.; Shalloway, D. Extending molecular dynamics time scales with milestoning: example of complex kinetics in a solvated peptide. *J Chem Phys*. 2007, 126, 145104-145104.
- [4] Haas, K.; Chu, J.W. Decomposition of energy and free energy changes by following the flow of work along reaction path. *J. Chem. Phys.* 2009, 131, 144105-144111.
- [5] Jnsson, H.; Mills, G.; Jacobsen, K.W. Nudged Elastic Band Method for Finding Minimum Energy Paths of Transitions, in Classical and Quantum Dynamics in Condensed Phase Simulations, Ed. B. J. Berne, G. Ciccotti and D. F. Coker, 385 (World Scientific, 1998).
- [6] E, W.; Ren, W.; Vanden-Eijnden, E. Simplified and improved string method for computing the minimum energy paths in barrier-crossing events. *J. Chem. Phys.* 2007, 126, 164103.
- [7] Dellago, C.; Bolhuis, P.G.; Geissler, P.L. Transition Path Sampling, *Adv. Chem. Phys.* 2002, 123, 1-84.
- [8] Cheng, X.; Wang, H.; Grant, B.; Sine, S.M.; McCammon, J.A. Targeted Molecular Dynamics Study of C-Loop Closure and Channel Gating in Nicotinic Receptors. 2006, 9, 134.
- [9] Elber, R. Long-timescale simulation methods. *Cur. Opin. in Str. Biol.* 2005, 15, 151-156.
- [10] Torrie, G. M.; Valleau, J. P. Nonphysical sampling distributions in Monte Carlo free-energy estimation: Umbrella sampling (1977) *J. Comput. Phys.* 23, 187
- [11] Tironi, I.G.; van Gunsteren, W.F. A molecular-dynamics simulation study of chloroform. *Mol. Phys.* 1994, 83, 381-403.
- [12] Meirovitch, H. Recent developments in methodologies for calculating the entropy and free energy of biological systems by computer simulation. *Current Opinion in Structural Biology*, 2007, 17, 181-186.

- [13] Mobley, D.L.; Chodera, J.D.; Dill, K.A. The combining and release method: obtaining correct binding free energies in the presence of protein conformational change. *Journal of Chemical Theory and Computation* 2007, 3, 1231-1235.
- [14] Mobley, D.L.; Klimovich, P.V. Perspective: Alchemical free energy calculations for drug discovery. *J. Chem. Phys.* 2012, 137, 230901-12.
- [15] Mobley, D.L.; Chodera, J.D.; Dill, K.A. On the use of orientational restraints and symmetry corrections in alchemical free energy calculations. *J. Chem. Phy.* 2006, 125, 084902.
- [16] Chipot, C.; Shell, M.S.; Pohorille, A. Introduction, in Chipot, C., Pohorille, A., editors. *Free Energy Calculations: Theory and Applications in Chemistry and Biology*. Springer Series in Chemical Physics, vol. 86. Berlin and Heidelberg: Springer; 2007, p. 132.
- [17] Jorgensen, W.L. The many roles of computation in drug discovery, *Science* 2004, 303, 18138.
- [18] Gilson, M.K.; Zhou, H.X. Calculation of protein-ligand binding affinities. *Annu Rev Biophys Biomol Struct.* (2007) 36, 21-42.
- [19] Handl, J.; Knowles, J.; Lovel, S.C. Artefacts and biases affecting the evaluation of scoring functions on decoy sets for protein structure prediction. *Bioinformatics*, 2009, 25, 1271-1279.
- [20] Dill, K.A.; H.S. Chan. From Levinthal to Pathways to Funnels: The "New View" of Protein Folding Kinetics. *Nature Structural Biology* 4, 10-19 (1997)
- [21] Dill, K.A.; Ozkan, S.B.; Shell, M.S.; Weikl, T.R. The protein folding problem. *Annual Review of Biophysics* (2008), 37, 289-316.
- [22] Anfinsen. C.B. Principles that Govern the Folding of Protein Chains. *Science* (1973) 181, 223-230.
- [23] Christ, C.D.; van Gunsteren, W.F. Enveloping distribution sampling: A method to calculate free energy differences from a single simulation, *J. Chem. Phys.* (2007), 126, 184110.

- [24] Ytreberg, F.; Zuckerman, D. Simple estimation of absolute free energies for biomolecules. *J. Chem. Phys.* 2006, 124, 104105.
- [25] Park, S.; Lau, A.; Roux, B. Computing conformational free energy by deactivated morphing. *J. Chem. Phys.* 2008, 129, 134102
- [26] Zheng, L.; Chen, M.; Yang, W. Random walk in orthogonal space to achieve efficient free-energy simulation of complex systems, *Proc. Natl. Acad. Sci.* 2008, 105 (51), 20227.
- [27] Shell, S.M. A replica-exchange approach to computing peptide conformational free energies. *Mol. Sim.* 2010, 7, 505-515.
- [28] Tyka, M.; Clarke, A.; Sessions, R. An Efficient, Path-Independent Method for Free-Energy Calculations. *J.Phys.Chem. B* 2006, 110, 17212-17220.
- [29] Cecchini, M., Krivov, S.V., Spicthy, M., Karplus, M. Calculation of free-energy differences by confinement simulations. Application to peptide conformers. *J. Phys. Chem. B.* 2009, 113, 9728-9740.
- [30] Ovchinnikov, V.; Cecchini, M.; Karplus, M. A Simplified Confinement Method for Calculating Absolute Free Energies and Free Energy and Entropy Differences. *J. Phys. Chem. B.* 2013, 117, 750-762.
- [31] Spicthy, M.; Cecchini, M.; Karplus, M. Conformational Free-Energy Difference of a Miniprotein from Nonequilibrium Simulations. *J. Phys. Chem. Lett.*, 2010, 1, 1922-1926.
- [32] Strajbl, M.; Sham, Y.Y.; Vill, J.; Chu, Z.-T.; Warshel, A. Calculations of Activation Entropies of Chemical Reactions in Solution. (2000) 104, 4578-4584.
- [33] Krivov, S.; Karplus, M. Hidden complexity of free energy surfaces for peptide (protein) folding *Proc. Natl. Acad. Sci. U.S.A.* 2004, 101, (41), 14766.
- [34] Alexander, P.A.; He, Y.; Chen, Y.; Orban, J. Bryan, P. The design and characterization of two proteins with 88% sequence identity but different structure and function. *Proc. Natl. Acad. Sci.* 2007, 104 (29), 11963-11968.

- [35] He, Y.; Chen, Y.; Alexander, P.A.; Orban, J. NMR structures of two designed proteins with high sequence identity but different fold and function. *Proc. Natl. Acad. Sci.* 2008, 105 (38), 14412-14417.
- [36] Alexander, P.A.; He, Y.; Chen, Y.; Orban, J. Bryan, P. A minimal sequence code for switching protein structure and function. *Proc. Natl. Acad. Sci.* 2009 , 106(50), 21149-21154.
- [37] Bryan, P.N.; Orban, J. Proteins that switch folds. *Curr Opin Struct Biol.* 2010, 20(4), 482-488.
- [38] He, Y.; Chen, Y.; Alexander, P.A.; Bryan, P.N.; Orban, J. Mutational tipping points for switching protein folds and functions. *Structure.* 2012, 20(2), 2 83-91.
- [39] Shortle, D. One sequence plus one mutation equals two folds. *Proc. Natl. Acad. Sci.* 2009, 106(50), 21011-21012.
- [40] Sheffler, W.; Baker, D. RosettaHoles: Rapid assessment of protein core packing for structure prediction, refinement, design, and validation. *Protein Sci ence.* 2009, 18(1), 229-239.
- [41] D.A. Case, T.A. Darden, T.E. Cheatham, III, C.L. Simmerling, J. Wang, R.E. Duke, R. Luo, R.C. Walker, W. Zhang, K.M. Merz, B. Roberts, S. Hayik, A. Roitberg, G. Seabra, J. Swails, A.W. Goetz, I. Kolossvry, K.F. Wong, F. Paesani, J. Vanicek, R.M. Wolf, J. Liu, X. Wu, S.R. Brozell, T. Steinbrecher, H. Gohlke, Q. Cai, X. Ye, J. Wang, M.-J. Hsieh, G. Cui, D.R. Roe, D.H. Mathews, M.G. Seetin, R. Salomon-Ferrer, C. Sagui, V. Babin, T. Luchko, S. Gusarov, A. Kovalenko, and P.A. Kollman (2012), AMBER 12, University of California, San Francisco.
- [42] Goetz, A.W.; Williamson, M.J.; Xu, D.; Poole, D.; Le Grand, S.; Walker, R.C. Routine microsecond molecular dynamics simulations with AMBER - Part I: Generalized Born. *J. Chem. Theory Comput.* 2012, 8(5) 1542.
- [43] MacCallum, J.; Perez, A.; Schnieders, MJ.; Hua, L.; Jacobson, M.P.; Dill, K.A. Assessment of protein structure refinement in CASP9. *Proteins*, 2011, 79, 74-90.

- [44] Kryshtafovych, A.; Fidelis, K; and Tramontano, A. Evaluation of model quality predictions in CASP9. *Proteins*, 2011, 79, 91106
- [45] Zemla, A. LGA: a method for finding 3D similarities in protein structures. *Nucleic Acids Res* 2003, 31, 33703374.
- [46] Perez, A.; Yang, Z.; Bahar, I.; Dill, K.A.; MacCallum, J.L.; FlexE: Using Elastic Network Models to Compare Models of Protein Structure. *J. Chem. Theory Comput.*, 2012, 8, 3985-3991.
- [47] Wang, Q.; Vantasin, K.; Xu, D.; Shang, Y. MUFOLD-WQA: A new selective consensus method for quality assessment in protein structure prediction. *Proteins*, 2011, 79: 185-195.
- [48] Case, D.A.; Cheatham, III, T.E.; Darden, T.; Gohlke, Luo, H.R.; Merz, Jr., K.M.; Onufriev, A; Simmerling, C.; Wang, B.; R. Woods, R. The Amber biomolecular simulation programs. *J. Computat. Chem.* (20005) 26, 1668-1688.
- [49] Hornak, V.; Abel, R.; Okur, A.; Strockbine, B.; Roitberg, A.; Simmerling, C. Comparison of multiple Amber force fields and development of improved protein backbone parameters. *Proteins.*, 2006, 65, 712-725.
- [50] Mongan, J.; Simmerling, C.; A. McCammon, J.; A. Case, D.; Onufriev, A. Generalized Born with a simple, robust molecular volume correction. *J. Chem. Theory Comput.*, 2006, 3, 156-169.

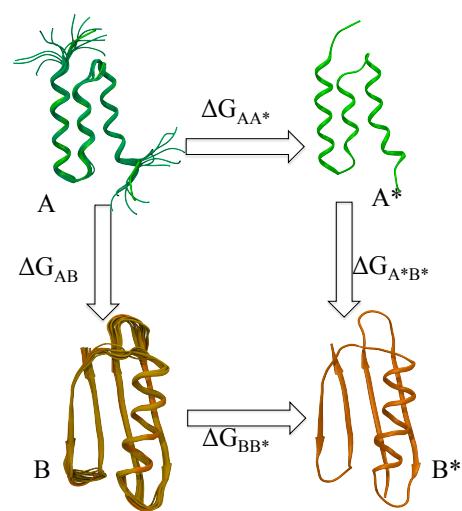
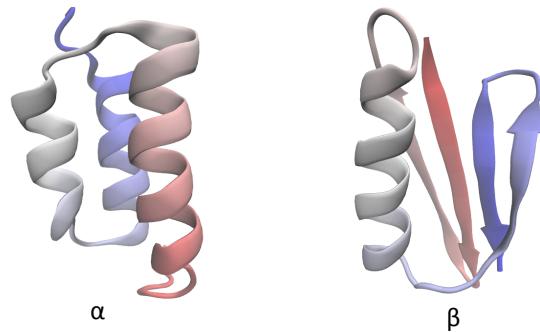


Figure 1: Graphical representation of the thermodynamic cycle involving confinement method.

12345678901234567890123456789012345678901234567890123456  
 2jws: TTYKLILNLQAKEEAIKELVDA**GIAEKYIKL**I<sup>A</sup>N<sup>K</sup>TVEGVWT**L**KDE**I**L<sup>T</sup>F<sup>F</sup>T<sup>V</sup>T<sup>E</sup> GA88  
 2jwu: TTYKLILNLQAKEEAIKELVDA**ATAEKYFKLY**<sup>A</sup>N<sup>K</sup>TVEGVWT**Y**KDET**K**T<sup>F</sup>T<sup>V</sup>T<sup>E</sup> GB88  
 2kdl: TTYKLILNLQAKEEAI**KELVDAGTAEKYIKL**I<sup>A</sup>N<sup>K</sup>TVEGVWT**L**KDE**I**KT<sup>F</sup>T<sup>V</sup>T<sup>E</sup> GA95  
 2kdm: TTYKLILNLQAKEEAI**KEAVDAGTAEKYFKL**I<sup>A</sup>N<sup>K</sup>TVEGVWT**Y**KDE**I**KT<sup>F</sup>T<sup>V</sup>T<sup>E</sup> GB95



Sequence	Experimental Fold		Calculated Fold		Calculated $\Delta G_{\beta} - \Delta G_{\alpha}$ (Kcal/Mol)
	$\alpha$	$\beta$	$\alpha$	$\beta$	
GA88	✓		✓		+3.94±0.51
GB88		✓		✓	-4.36±0.46
GA95	✓		✓		+3.48±0.47
GB95		✓		✓	-5.01±0.49

Figure 2: The confinement method correctly predicts the structural preferences of four chameleon sequences. The top part of the figure represent four sequences used in this study along with the protein data bank identifier. The experimentally observed fold, computationally predicted fold, and difference free energy between the two folds is reported for each sequence.

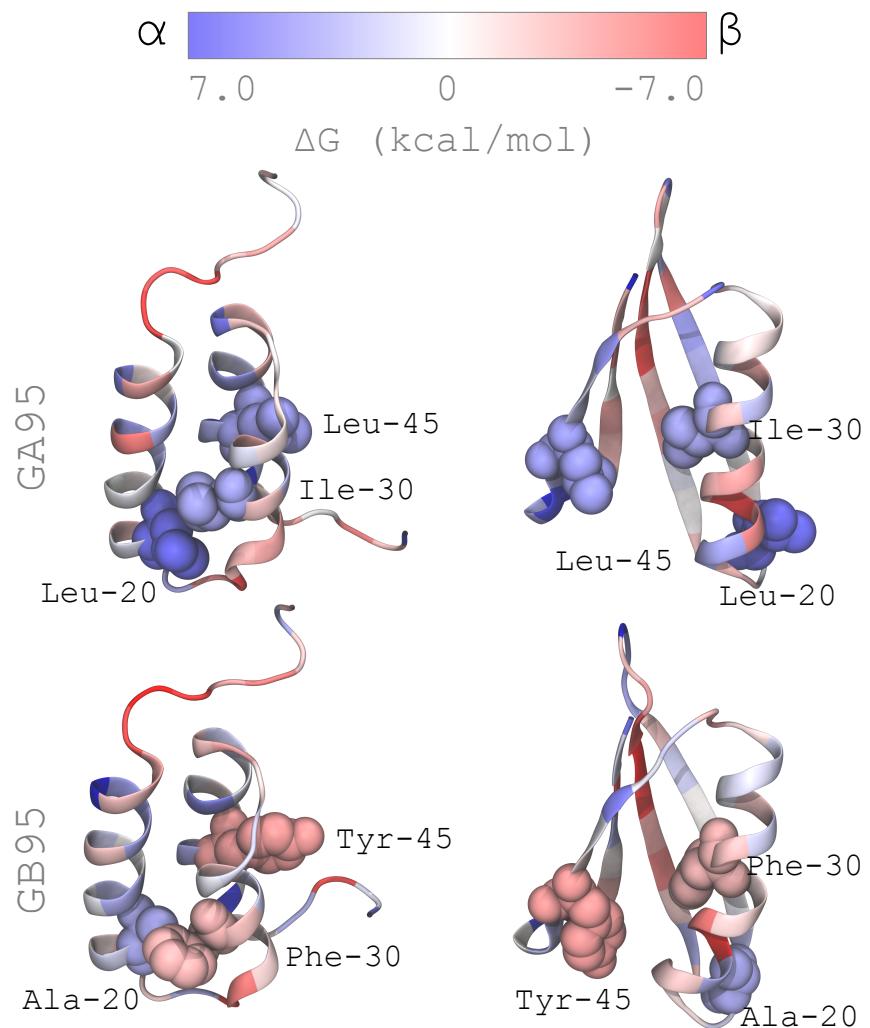


Figure 3: Per-residue difference free energies for the  $\alpha$  and  $\beta$  conformations of the GA95 and GB95 sequences. Residues colored in blue favor the  $\alpha$  structure, residues in red favor the  $\beta$  structure, and white residues have no preference.

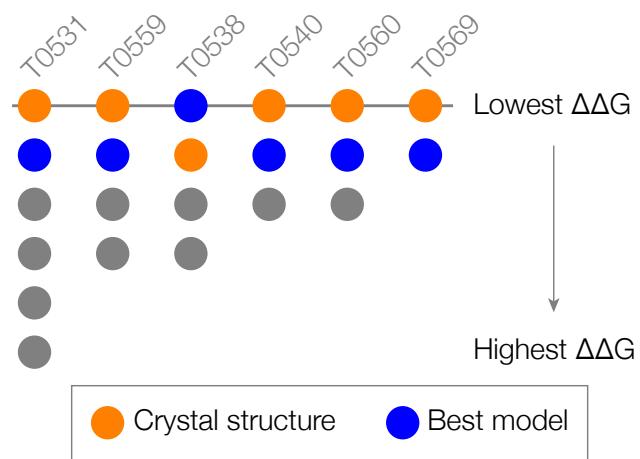


Figure 4: The confinement method is usually able to identify the native structure and the best model (the model with the highest GDT-TS) from a set of decoys.

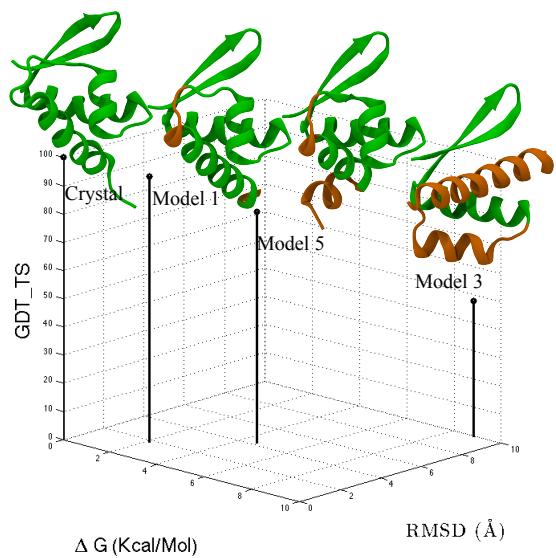


Figure 5: The confinement method can correctly rank the native structure and three predictions produced by a single method.

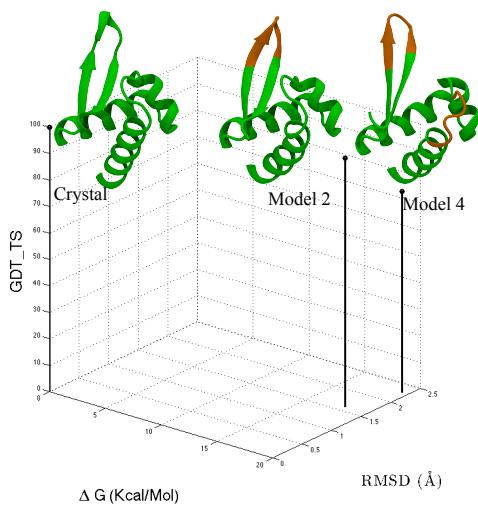


Figure 6: Native structure and two models (from group “Splicer”) for target T0560.

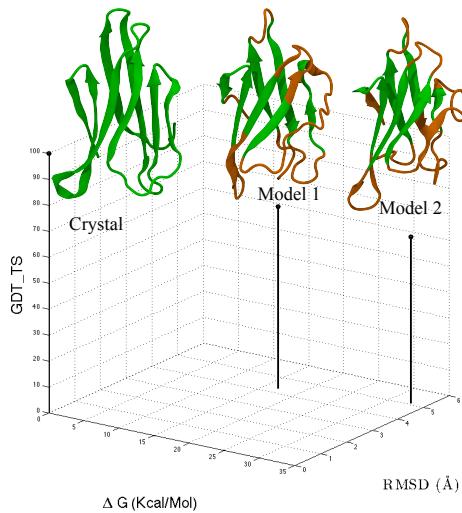


Figure 7: The confinement method can correctly rank the native structure and two models submitted by different prediction groups.

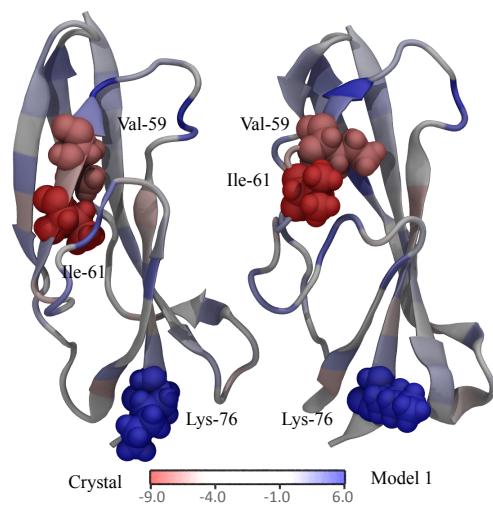


Figure 8: Per-residue difference free energy between the experimental NMR structure and the best prediction for CASP target T0569. The amino acid residues that are colored in deep red and deep blue stabilizes the NMR structure and the prediction, respectively; the residues with light blue color does not have a strong preference.

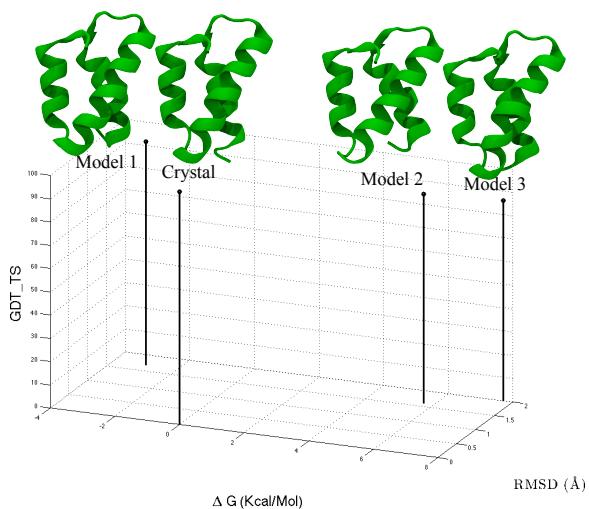


Figure 9: The native and three model structure of engineered protein from Asr4154 protein (PDB ID: 2L09 and CASP code:T0538). The model 1,2 and 3 are from the group PconsR, Shell and FOLDIT respectively.

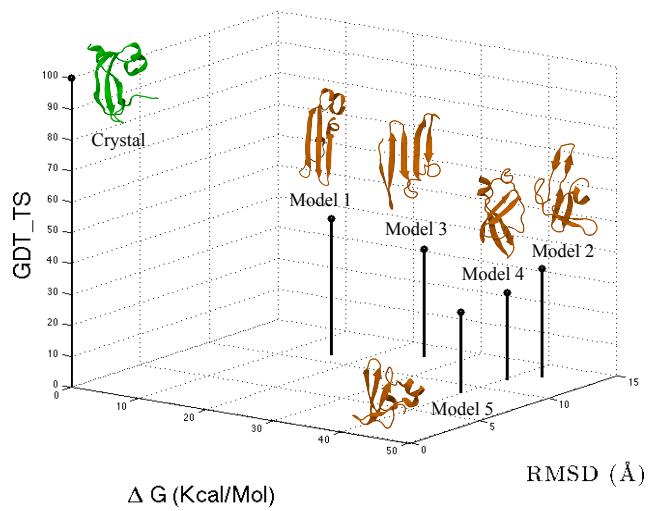


Figure 10: The native structure and 5 models of extracellular domain of the jumping translocation breakpoint protein (pdb id: 2KJX and the CASP code: T0531).

Table 1: The confinement method assigns a more favorable free energy to the experimentally determined structure than to computer-generated predictions. For each target, we examined as many as five predictions submitted by CASP participants. We report the free energy difference between the most favorable decoy and the experimentally determined structure. Positive  $\Delta\Delta G$  values indicate that the experimental structure is predicted to be more favorable than any of the decoys.

CASP Target	PDB Identifier	$\Delta\Delta G = G_{best\ decoy} - G_{native}$ (kcal/mol)
T0531	2KJX	$11.15 \pm 0.70$
T0538	2L09	$-3.00 \pm 0.47$
T0540	3MX7	$16.94 \pm 0.49$
T0559	2L01	$2.10 \pm 0.24$
T0560	2L02	$18.00 \pm 0.49$
T0569	2KWY	$20.01 \pm 0.69$