

Frustration in the energy landscapes of multidomain protein misfolding

Weihua Zheng^{a,b}, Nicholas P. Schafer^{b,c}, and Peter G. Wolynes^{a,b,c,1}

Departments of ^aChemistry and ^cPhysics, and ^bCenter for Theoretical Biological Physics, Rice University, Houston, TX 77005

Contributed by Peter G. Wolynes, December 18, 2012 (sent for review November 29, 2012)

Frustration from strong interdomain interactions can make misfolding a more severe problem in multidomain proteins than in single-domain proteins. On the basis of bioinformatic surveys, it has been suggested that lowering the sequence identity between neighboring domains is one of nature's solutions to the multidomain misfolding problem. We investigate folding of multidomain proteins using the associative-memory, water-mediated, structure and energy model (AWSEM), a predictive coarse-grained protein force field. We find that reducing sequence identity not only decreases the formation of domain-swapped contacts but also decreases the formation of strong self-recognition contacts between β -strands with high hydrophobic content. The ensembles of misfolded structures that result from forming these amyloid-like interactions are energetically disfavored compared with the native state, but entropically favored. Therefore, these ensembles are more stable than the native ensemble under denaturing conditions, such as high temperature. Domain-swapped contacts compete with self-recognition contacts in forming various trapped states, and point mutations can shift the balance between the two types of interaction. We predict that multidomain proteins that lack these specific strong interdomain interactions should fold reliably.

aggregation | funnel

Protein misfolding and productive protein folding bear a yin-yang relationship in the energy landscape theory of biomolecular self-organization (1). Only by comparing the strengths of the forces leading to proper structure to those that might, by chance, stabilize alternative structure can we quantitatively understand how proteins kinetically access their thermodynamically stable ordered states (1). In vivo and at low concentrations in vitro, unfolded small proteins avoid kinetic traps and generally find their way easily to their native state. Nevertheless, diseases caused by the misfolding of several specific proteins plague mankind (2, 3). Despite much effort, the patterns of interactions that allow pathological misfolding remain incompletely understood. Known pathological misfolding entails aggregation of specific proteins and thus the interactions of protein molecules with other copies of themselves. Energy landscape theory provides one natural explanation of this specificity in misfolding through the funneled nature of the monomeric protein energy landscape: Native-like interactions between different protein molecules like those found within a single protein are stronger than alternate nonnative interactions in the same molecule or interactions between peptide sequences chosen at random in the two molecules. Because of this intrinsic self-stickiness of foldable molecules, runaway domain swapping, in which native-like interactions are made between different copies of the same protein, provides a natural mechanism for aggregation (4–7). Indeed, transient protein aggregation during refolding at moderately high concentration does appear to be universal (8). Nevertheless this aggregation usually resolves itself eventually as the system comes to equilibrium, thus arguing that something more may be involved in natural pathological misfolding that seems permanent. One attractive idea is that there is an alternate “amyloid funnel” (9), which a protein may enter if the molecule has enough time to find it before completing its native folding. A funnel to the

amyloid state has been thought to possibly be universal, because under appropriate denaturing conditions, it seems, even the most innocuous proteins can form amyloids (10). Alternatively, like the funnel for formation of a native structure, the amyloid funnel may be encoded in sequence signals (11, 12). We have been led to address these questions about the misfolding energy landscape in our effort to model a series of insightful experimental investigations on multidomain proteins (13, 14). Multidomain proteins are much more susceptible to misfolding than single-domain proteins because of the effective high local concentration of peptide binding partners. Experiments on artificial constructs in which related protein domains are fused together indicate that high sequence identity of the domains favors aggregation and that the initial interactions of the fused domains are critical to the aggregation process (13). At the same time, bioinformatic studies indeed show that neighboring domains in natural multidomain proteins have lower than expected sequence identity (15). These observations point to the importance of domain swapping as predicted by the minimal frustration principle from energy landscape theory (16). The present computational studies show that this is only part of the story, however. Indeed in our simulations additional sequence signals that allow some peptide fragments to recognize copies of themselves greatly increase the tendency to misfold. In silico, single mutations in these fragments can significantly reduce misfolding of a multidomain construct. It turns out that these sequence signals do not yield a globally unique structure but are able to act on rather high-entropy ensembles. They do not affect the ordered native state ensembles. These sequence signals allow alternate structures of self-recognizing peptide fragments to achieve minimally frustrated configurations that are locally competitive with the final native structure. Making these alternate structures, however, frustrates the formation of further tertiary structure in the protein so these misfolded states are globally higher in energy than the native state: Stable misfolding is encouraged by the high residual entropy of an ensemble of structures that can take advantage of a locally strong interaction. Misfolding thus occurs optimally under intermediate denaturing conditions. In this scenario, then, pathological aggregation remains consistent with a largely minimally frustrated and funneled energy landscape of each monomer. The residual frustration that encourages misfolding is not evident at the residue pair level but involves alternate pairings of hexamer fragments, allowing self-recognition in amyloid-like assemblies.

The computational approach we take is based on transferable energy functions that have been optimized to predict protein tertiary structures of monomers using simulated annealing (17). To do this annealing we leverage the associative-memory, water-mediated, structure and energy model (AWSEM)-MD software

Author contributions: W.Z., N.P.S., and P.G.W. designed research; W.Z. and N.P.S. performed research; N.P.S. contributed new reagents/analytic tools; W.Z. and N.P.S. analyzed data; and W.Z., N.P.S., and P.G.W. wrote the paper.

The authors declare no conflict of interest.

¹To whom correspondence should be addressed. E-mail: pwolynes@rice.edu.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1222130110/-DCSupplemental.

package that has been shown to accurately predict monomeric (18) and properly oligomeric protein structures (19). Although the energy function is bioinformatically based, it appears to have many elements of biophysical realism. Here we show the same energy function encodes also the local signals for misfolding, even though information on such misfolding was not used in training the algorithm. In addition to simulating the rapid annealing of several fused dimer constructs, we compute multidimensional free-energy profiles under varying thermodynamic conditions to understand the entropy/energy interplay on the landscape. We also quantify the statistics of the conformational states formed with both native structures and elements of misfolded structure. The self-recognition ensemble is energetically less stable than the native ensemble because the abnormally strong self-recognition interaction is highly local. The same energy function used for simulated annealing allows us to rapidly scan any sequence for fragments that are likely to participate in such amyloid-like structure formation. The results of this scan are consistent with other bioinformatic tools for identifying amyloidogenic sequences (11, 20–24).

Results and Discussion

Simulated Annealing of Fused Multidomain Proteins: Misfolding Correlates with the Sequence Identity Between Domains. A couple of recent and insightful experimental studies of misfolding and aggregation have focused on the Ig domains of the vertebrate muscle protein titin (13, 14). One study showed that the rate of aggregation of the fused dimers is correlated with the sequence identity between the domains (13). Our simulation investigation parallels this experimental work. We investigated fused dimers in which the 27th Ig domain of human cardiac titin [TI I27; Protein Data Bank (PDB) ID 1TIU] is the first domain. Five proteins with varying sequence identity to I27 were chosen for the second domain, which is connected to the first domain via a four-residue glycine linker. In addition, we also simulated an SH3-SH3 fused dimer. Starting from totally extended conformations, 40 independent simulated annealing simulations were run for each fused protein. Simulated annealing searches for energetically stable structures by gradually reducing the temperature from slightly above the folding temperature, at which the efficiency of searching the conformational space is considered optimal, to the native temperature at which the native structure is stable. Our simulated annealing protocol (*SI Text*) is not an attempt to reproduce faithfully the experimental conditions (13, 14), but rather is used as a way of quickly searching the conformational landscape for states that could trap the protein and thereby inhibit or slow productive folding. At the end of a simulated annealing run, the protein typically adopts either the native conformation or a compact misfolded conformation such as the self-recognition state in the case of the I27-I27 fused dimer. Using the final structures from these simulations, we calculated the fraction of misfolded domains. Fig. 1 shows the fraction of misfolded domains vs. the sequence identity between the two domains. A domain is considered folded when the fraction of native contacts within the domain $Q_{\text{domain}}^i > 0.4$, $i = 1, 2$. As the sequence identity increases, the fraction of misfolded domains that result from simulated annealing increases, consistent with the observation that the sequence identity between neighboring domains in natural multidomain proteins is lower than would be expected by chance (13). However, what is the nature of these misfolded ensembles? What types of interactions are responsible for the misfolding?

Self-Recognition Contacts and Domain-Swapped Contacts in the Misfolded Structures. Interchain interactions are important for folding and misfolding when the protein concentration is high. In the case of fused dimer systems, the local concentration is always high—one domain has ample opportunity to interact with its covalently linked neighboring domain even when the concentration

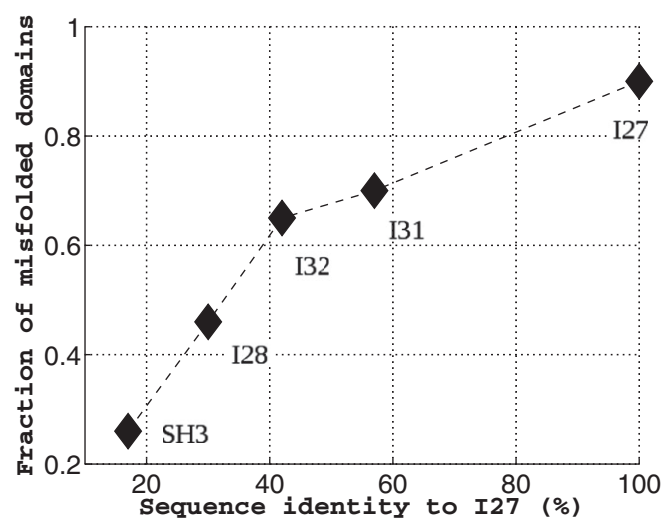


Fig. 1. The fraction of misfolded domains found upon simulated annealing is plotted against the sequence identity between the two domains in the fused construct. Titin I27 was fused with five different proteins via a four-residue GLY linker to form two-domain proteins, as indicated in the plot. Forty annealing simulations were run for each fused protein from totally extended conformations. All proteins were studied with the same annealing schedule. A domain is considered folded when the fraction of native contacts within the domain $Q_{\text{domain}}^i > 0.4$, $i = 1, 2$ at the end of all of the simulations. As the sequence identity increases, the fraction of misfolded domains increases accordingly.

of molecules is low. Wright et al. have shown that, in the case of fused dimers of the Ig domains of titin, this initial interaction between the fused domains is the important step for determining aggregation rates (13). In these experiments, the aggregation rate was insensitive to the number of copies that were fused together for $n \geq 2$. However, the nature of the misfolded structures remains unclear. Several a priori extreme possibilities exist—either specific interactions could drive misfolding or the misfolded structures might appear completely disordered and random. For evolved proteins, which satisfy the principle of minimal frustration, domain-swapped interactions are the most obvious candidate for specific interactions that drive misfolding. The counterparts of domain-swapped contacts in the monomer are native contacts that are in general stronger than other contacts and allow the protein to fold quickly and reliably. These same strong interactions can also drive oligomerization via domain swapping with nearby domains. Another candidate for misfolding, formation of self-recognition contacts, has been identified in studies of amyloid fibrils. Microcrystals of fibril-forming proteins reveal a “steric zipper” structure with two self-complementary β -sheets that form a spine of an amyloid fibril (25). These self-recognition contacts between two amyloidogenic segments in different molecules can be extremely strong but, unlike domain-swapping interactions, have no exact counterpart in the native structure and therefore can be involved only in misfolding. These segments are rich in hydrophobic amino acids, and bioinformatic studies indicate that evolution suppresses long stretches of hydrophobic amino acids (26). However, evolution has not completely eliminated them—these “amyloidogenic” segments appear to exist at a frequency of about one per protein (27) and are almost always buried in the natively folded structure. When the concentration of a particular protein is low, as is commonly the case in vivo, self-recognition of two buried segments is unlikely. Therefore, it is unsurprising that not every protein that has an amyloidogenic segment is involved in forming pathological fibrils or aggregates in vivo. Nonetheless, it is of great interest to study the role of these interactions in misfolding.

a strand of seven residues from monomer A with the same seven residues of monomer B. A calculation using the AWSEM-Amylometer (*SI Text*) identifies that these seven residues include a hexamer 56–61 (HILILH) that has the strongest self-recognition interaction among all hexamers in the sequence of I27. Furthermore, this self-recognition interaction is stronger than any possible native or other nonnative hexamer pair interaction in I27-I27 and is lower in energy than the threshold for amyloidogenicity as determined by the AWSEM-Amylometer.

Misfolded Ensembles Stabilized by Self-Recognition Interactions Are Entropically More Favored Than the Native. Fig. 3 shows the energy and free-energy landscapes of the I27-I27 fused dimer along two reaction coordinates: a native reaction coordinate, the fraction of native contacts Q and a nonnative reaction coordinate, the sum of the number of self-recognition contacts N_{self} and the number of swapped contacts N_{swap} . The misfolded state I is energetically less stable than the native state N but entropically more favored. As shown in Fig. S1, below the folding temperature the native ensemble is the most populated state. As temperature increases, the misfolded ensemble becomes more and more highly populated. This type of metastable ensemble acts as a kinetic trap even below the folding temperature. When the fused dimer construct becomes trapped in this metastable state, large stretches of the structure are disordered and parts of the protein that are normally buried become exposed to solvent. This type of misfolding event could lead to aggregation when many fused dimers are present and is consistent with the experiments of Wright et al. (13).

Point Mutations and Their Effects in I27-I27 and SH3-SH3. Removal and addition of strong self-recognition interactions by point mutations can change the degree of misfolding significantly. The foldability of the fused dimers correlates well with the sequence identity between the two domains as shown in Fig. 1. As the sequence identity of neighboring domains decreases, the fused protein folds correctly more often. In simulated annealing of I27-I27, strong self-recognition contacts often result in misfolding.

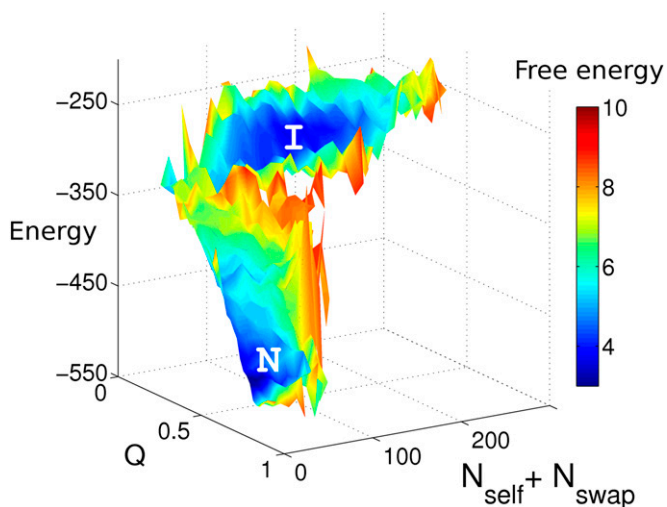


Fig. 3. Energy and free-energy surfaces for I27-I27 at its folding temperature. N_{self} and N_{swap} are the number of self-recognition contacts and the number of domain-swapped contacts, respectively. The trapped states I have higher energies than the native states N, as shown in the z axis, but have similar free energies to those of the native states, as shown by the color coding of the free energy, with scale indicated in the side bar. We see that the ensemble I states are entropically favored. As temperature increases, the intermediate ensemble will become more stable than the native ensemble.

Can we mutate the sequence to eliminate these strong interdomain interactions so that the protein folds more often? We constructed a double mutant of I27, V11D and I59E. The resulting I27*⁻I27*^{*} has 100% sequence identity between the two domains, but two of the strongest interdomain hexamer pairs that were found most frequently in previous misfolded structures in I27-I27 simulated annealing have been greatly weakened. The AWSEM-Amylometer was used to find the appropriate mutations that would reduce the strength of these self-recognition interactions the most. Despite having 100% sequence identity, I27*⁻I27*^{*} folds three times more often than the I27-I27 system and folds as often as fused I27-I31 (57% sequence identity).

For SH3-SH3, which folds well, we tested to see whether conversely a single-point mutation that introduces a strongly self-recognizing hexamer can induce misfolding in the simulated annealing. Again, the AWSEM-Amylometer was used to suggest an appropriate mutation. After introducing the E27I mutation, the fraction of misfolded domains goes up from 18% to 79%, using the identical simulation protocol that was used before. These two mutation studies suggest that one of the important effects of lowering the sequence identity between neighboring domains is to decrease the probability of strongly self-recognizing segments occurring, which in turn lowers the chance of misfolding via interdomain interactions.

Competing Roles of Self-Recognition Contacts and Swapped Contacts.

As shown in Fig. 3, a combination of self-recognition contacts and swapped contacts contributes to the stabilization of misfolded ensembles. Some self-recognition contacts are particularly strong. Swapped contacts are also strong as would be anticipated by the principle of minimal frustration. By making point mutations, we can change the relative strength of both interactions and study their competing roles in forming various misfolded structures. As shown in Fig. 4, the I27-I27 fused dimer, where strong self-recognition interactions are present, favors misfolded structures with these interactions satisfied. Forming these contacts in turn inhibits misfolded structures with many swapped contacts from forming. When point mutations were made to reduce the stability of these hexamer pairs, the structures with a significant number of swapped contacts appear more often. The converse is true for SH3-SH3. The point mutation E27I introduces a strong hexamer pair, promoting the formation of misfolded structures with the hexamer pair forming instead of forming domain-swapped structures. Both types of contacts can be strong, and depending on the specific sequence and simulation conditions, one or both can contribute to misfolding.

The single-molecule experiments of Borgia et al. (14) reveal some misfolding of I27-I27 while refolding under native conditions. Under these conditions only a small amount of misfolded proteins is formed. Using a symmetrized Go model Borgia et al. were able to propose several candidate misfolded (domain-swapped) structures and found that the FRET distances that they measured were consistent with these structures, although they were not able to distinguish between them. Interestingly, one of the domain-swapped structures that they observed in simulation was also found in our simulated annealing simulations of I27*⁻I27*^{*}, but was not found by us for I27-I27. In our simulations of I27-I27, the self-recognition hexamer pairing is so strong that it suppresses the formation of a fully domain-swapped structure. In I27*⁻I27*^{*}, the weakening of the self-recognition interactions allows domain-swapped interactions to predominate. To compare our structures with the experimental results of Borgia et al., we took the self-recognition ensemble of misfolded structures from the I27-I27 simulations and measured the distance between the residues that were labeled in the experiments, residues 3 and 83. The distribution of distances peaks around 2.3 nm, as shown in Fig. S2, very close to 2.0 nm, the distance in the domain-swapped structure proposed in their paper. For further experiments to

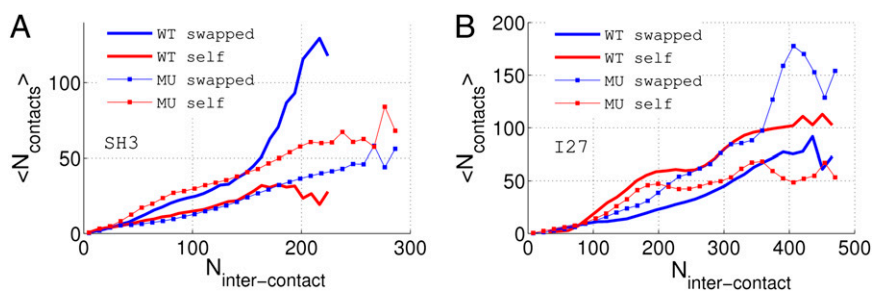


Fig. 4. Competing roles of the swapped contacts and self-recognition contacts in the “wild-type” (WT) fused dimer I27-I27, SH3-SH3, and their mutants (MT) I27*-I27* and SH3*-SH3*. Forty annealing simulations were carried out for all dimers. $\langle N_{\text{contacts}} \rangle$ is the number of contacts averaged over all 40 simulations. $N_{\text{intercontact}}$ is the number of interdomain contacts. The interdomain contacts include the swapped contacts, self-recognition contacts, and other types of interdomain contacts. (A) Single-point mutation introduces a pair of strong self-recognition contacts in SH3*-SH3*, suppressing the formation of swapped contacts. Significant misfolding with formation of self-recognition contacts occurs after the mutation. (B) Two point mutations eliminated two pairs of the strongest self-recognition contacts in I27; therefore the swapped contacts play a more dominant role.

clarify the detail of the misfolded structure, we suggest that two FRET labels be put near the strongly self-recognizing hexamer of each monomer. If the self-recognition contacts are in fact formed as predicted by the rapid simulated annealing of the present model and also remain stable under experimental conditions, the transfer efficiency should be high and strongly peaked for the newly labeled misfolded construct.

Folding and Misfolding from Single Domain to Multidomain: Implications for Aggregation. In single-domain protein folding, the native interactions are in competition with nonnative intrachain interactions. Evolution has made these nonnative intrachain interactions weaker than the native ones so that proteins fold correctly on biological timescales. The misfolding problem gets more severe in the case of multidomain proteins or proteins at high concentration because of the possibility of forming strong interdomain contacts. Nearly all protein sequences have at least one amyloidogenic segment (27). Due to the strong interactions between them, multiple copies of these segments in proximity could cause the protein or proteins to adopt transiently stable trapped conformations that trigger further aggregation. For I27, as shown in Fig. 5A, inclusion of interdomain interactions in the AWSEM-Amylometer calculation reveals a specific self-recognizing hexamer interaction that is amyloidogenic according to our empirically determined threshold. It is important to remember that this interaction is simply not possible in the case of the monomer itself because only a single copy of the hexamer exists in a monomer. When simulated annealing runs were performed on the I27-I27 fused dimer, it was indeed this particular interaction that was the dominant cause for misfolding, as discussed previously. On the other hand, SH3-SH3 folds as well as its monomer because all of the hexamers in SH3 are only weakly self-recognizing. Our folding and misfolding results on the I27 and SH3 fused dimers are consistent with the experimental aggregation studies that show that the I27-I27 fused dimer aggregates but SH3-SH3 does not (13).

Self-recognizing segments tend to be short, typically between five and seven residues long, and so the interactions that stabilize the resulting misfolded structures are highly localized in sequence. This short sequence length is about that of a Kuhn statistical segment of a polypeptide (29) so the entropy loss in such self-recognition is comparable to that of bringing a single pair of residues together. This allows self-recognition to occur locally. Therefore, as shown in Fig. 5C, even though the local self interactions between the hexamers are much stronger than the average native interactions, the self-recognition ensemble is energetically less stable than the native ensemble. Swapped contacts are not as strong as the strongest self-recognition contacts but are present diffusely throughout the sequence. Forming

domain-swapped structures requires the cooperation of contacts throughout the sequence but results in structures that are closer in energy to the natively folded structure. Both of these types of contacts can lead to oligomerization, but the resulting misfolded structures have different thermodynamic contributions to their stability. To what extent these two types of interaction contribute to the different stages of pathological aggregate formation in vivo remains an open question, but the coincidence of detecting amyloid diseases and occurrence of fevers in patients is intriguing (30).

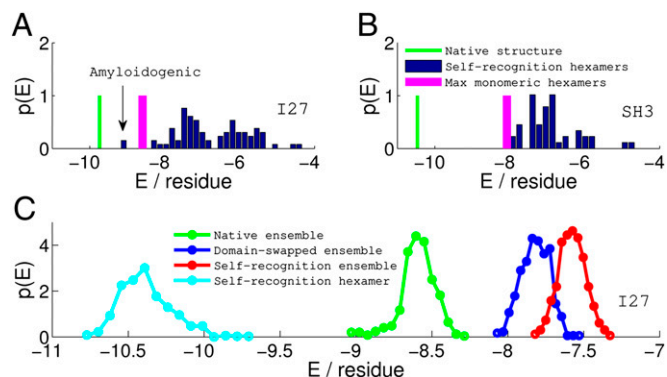


Fig. 5. Comparisons of stability (energy per residue) among various zero-temperature structures and ensembles of thermally sampled structures of I27 (A and C) and SH3 (B). The stability for the native monomeric structure (green vertical line) is calculated from the AWSEM energy function. The strongest nonnative hexamer pairing possible in the monomer (magenta bar) is significantly less stable than the native structure, indicating that misfolding by inappropriate pairing of strands will be unlikely during the folding of the monomer for both I27 and SH3. In A and B, the blue bars represent the distribution of the stability of all of the self-recognition hexamer pairs, calculated from the AWSEM-Amylometer (SI Text). If the stability of the strongest self hexamer pair is competitive with the native structure, as in the case of I27-I27, the particular self pair becomes responsible for the misfolding of the fused protein in our simulation and potentially, would trigger further aggregation in solution. For SH3-SH3, B predicts that fused protein should fold as well as the monomer in the simulation, because all self hexamer pairings are weaker than the most stable nonnative hexamer pairing in the monomer. In C, the stability distributions of various ensembles of structures collected from the simulations of I27-I27 are shown. The native ensemble (green) is energetically more stable than both the domain-swapped ensemble (blue) and the self-recognition ensemble (red). Nevertheless, the local interactions between the self-pairing hexamers from the self-recognition ensembles, shown in cyan, are even stronger than typical energies in the native folded ensemble.

