

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/5302866>

# Microscopic Reversibility of Protein Folding in Molecular Dynamics Simulations of the Engrailed Homeodomain †

ARTICLE *in* BIOCHEMISTRY · AUGUST 2008

Impact Factor: 3.02 · DOI: 10.1021/bi800118b · Source: PubMed

---

CITATIONS

35

---

READS

26

3 AUTHORS, INCLUDING:



David A. C. Beck

University of Washington Seattle

51 PUBLICATIONS 987 CITATIONS

SEE PROFILE



Valerie Daggett

University of Washington Seattle

223 PUBLICATIONS 12,613 CITATIONS

SEE PROFILE

Published in final edited form as:

Biochemistry. 2008 July 8; 47(27): 7079–7089. doi:10.1021/bi800118b.

## Microscopic reversibility of protein folding in molecular dynamics simulations of the engrailed homeodomain†

Michelle E. McCully<sup>‡</sup>, David A.C. Beck<sup>§</sup>, and Valerie Daggett<sup>\*,‡,§</sup>

<sup>‡</sup>Biomolecular Structure and Design Program, University of Washington, Box 355013, Seattle, Washington 98195-5013

<sup>§</sup>Department of Bioengineering, University of Washington, Box 355013, Seattle, Washington 98195-5013

### Abstract

The principle of microscopic reversibility states that at equilibrium the number of molecules entering a state by a given path must equal those exiting the state via the same path under identical conditions, or in structural terms, that the conformations along the two pathways are the same. There has been some indirect evidence indicating that protein folding is such a process, but there have been few conclusive findings. In this study, we performed molecular dynamics simulations of an ultra-fast unfolding and folding protein at its melting temperature in order to observe, on an atom-by-atom basis, the pathways the protein followed as it unfolded and folded within a continuous trajectory. In a total of 0.67  $\mu$ s of simulation in water, we found 6 transient denaturing events near the melting temperature (323K and 330K) and an additional refolding event following a previously identified unfolding event at high-temperature (373K). In each case unfolding and refolding transition state ensembles were identified, and they agreed well with experiment based on comparison of S- and  $\Phi$ -values. Based on several structural properties, these 13 transition state ensembles agreed very well with each other and with 4 previously identified transition states from high-temperature denaturing simulations. Thus, not only were the unfolding and refolding transition states part of the same ensemble, but in five of the seven cases, the pathway the protein took as it unfolded was nearly identical to the subsequent refolding pathway. These events provide compelling evidence that protein folding is a microscopically reversible process. In the other two cases, the folding and unfolding transition states were remarkably similar to each other, but the paths deviated.

In 1925, Richard C. Tolman coined the term “microscopic reversibility” in reference to chemical reactions:

In recent years increasing use has been made of a new postulate which perhaps cannot yet be stated in its final form, but which requires in a general way in the case of a system in thermodynamic equilibrium not only that the total number of molecules leaving a given state in unit time shall on the average equal the number arriving in that state in unit time, but also that the number leaving by any particular path shall on the average be equal to the number arriving by the reverse of that particular path, thus excluding any cyclical maintenance of the equilibrium state. The writer has ventured to name this postulate *the principle of microscopic reversibility* (1).

This description was recast into structural terms in 1967 by Frank H. Westheimer (2) and was adopted by the IUPAC in 1999:

<sup>†</sup>This research was supported by the National Institutes of Health Grant GM50789 (to V.D.).

\* To whom correspondence should be addressed. daggett@u.washington.edu. Phone: (206) 685-7420. Fax: (206) 685-3300.

In the case of SN2 reactions at tetrahedral centers implying a formation of the trigonal bipyramid transition state (or intermediate) structure, the original formulation of the principle was extended in the following way: if a molecule or reactant enters a trigonal bipyramid at an apical position, this (or another) molecule or reactant must likewise leave the trigonal bipyramid from an apical position (3).

The hypothesis that protein folding may, like chemical reactions, be microscopically reversible has since been offered. If this hypothesis is true, one would expect to observe identical transition states for folding and unfolding, and major events on the folding pathway would occur in reverse order in the unfolding pathway. Supporting evidence has been presented by Jackson *et al.* who showed through  $\Phi$ -value analysis on chymotrypsin inhibitor 2 (CI2)<sup>1</sup> that folding and unfolding transition states are the same, suggesting that the pathways are also the same (4). Additionally, molecular dynamics (MD) generated unfolding transition states (TS) of CI2 are in quantitative agreement with experimental data collected for both unfolding and refolding (5,6). Furthermore, the unfolding and direct refolding pathways of CI2 were shown to be the same in a single continuous MD trajectory (7).

The latter MD study was done at the melting temperature ( $T_m$ ) of the protein, at which point the folding and unfolding rates are equal and  $\Delta G = 0$ . For these reasons, exchange between the folded and unfolded states is dependent on the energy barrier,  $\Delta G^\ddagger \approx 2.3$  kcal/mol (8), and unfolding and refolding may occur in a single trajectory on time scales tractable by MD. The protein, CI2, passed through three different states, native (N), nearly native (N'), and denatured (D), then returned to N' over a time period of about 60 ns. When moving from N' to D and back, unfolding and refolding transition states were identified, and they were the same. The C $\alpha$  root mean square deviation (RMSD) to the native structure and internal contacts were analyzed to differentiate between the three different states. Day and Daggett (7) defined N', an alternate, stable state for CI2 at elevated temperature. The protein passed through N' before it moved through its TS to D. N' was characterized by many near-native interactions but elongated contact distances. In particular, Trp 5, a fluorescence unfolding probe, was buried in both N and N', but not in D, as would be expected if both N and N' were native but D was not. D had a disrupted hydrophobic core and loss of secondary structure. This CI2 study showed direct unfolding and refolding in a single continuous trajectory by the same structural pathways for the first time. Consequently, this behavior needs to be demonstrated in another system to ensure it is reproducible, which we describe here.

The engrailed homeodomain (En-HD) of *Drosophila melanogaster* is a 61-residue three-helix bundle. It is ultra-fast folding ( $k_F = 37,500$  s<sup>-1</sup> at 25°C and 51,000 s<sup>-1</sup> around 42°C) and unfolding ( $k_U = 1,100$  s<sup>-1</sup> at 25°C and 205,000 s<sup>-1</sup> at 63°C), and its folding and unfolding pathways have been extensively characterized through combined experimental and MD studies (9-12). Folding for En-HD follows the framework model involving the docking of HI (residues 10-22), HII (28-38), and HIII (42-55) (13). These properties make En-HD especially well-suited for MD folding studies.

In this study, we performed MD simulations of En-HD near its  $T_m$  (52°C = 325K (12)) to compare unfolding and refolding under identical conditions. We analyzed 5 simulations, 3 at 323K and 2 at 330K. We compared them to 4 previously described thermal denaturation simulations, 2 each at 373K and 498K, and one native simulation at 298K. The first of the 373K simulations was found to contain a region of particular interest that had not been previously reported, so that simulation was also analyzed in detail. We identified and characterized 3 different states populated during unfolding and refolding: N, N', and D. We

<sup>1</sup>Abbreviations: CI2, chymotrypsin inhibitor 2; MD, molecular dynamics; TS, transition state;  $T_m$ , melting temperature;  $\Delta G^\ddagger$ , energy barrier / activation energy; N, native; N', nearly native; D, denatured; RMSD, root mean square deviation; En-HD, engrailed homeodomain; TSE, transition state ensemble; *ilmm*, *in lucem* molecular mechanics; MDS, multidimensional scaling.

also found 6 transient denaturing events in which En-HD partially unfolded and refolded in 3 of the 5  $T_m$  simulations; from this 12 unfolding and refolding transition state ensembles (TSE) were identified. En-HD unfolded in the 373K/1 simulation, and a TSE in agreement with experiment was reported previously (9,12,13). Further investigation of this simulation showed that En-HD later refolded, so we also describe this high temperature refolding TSE. These 13 TSEs agree well with the 4 previously identified unfolding TSEs from the 4 high-temperature simulations. Besides defining TSEs, we analyzed the entire pathway En-HD followed as it unfolded and refolded. Five of the 7 refolding pathways were nearly identical to the unfolding pathways that preceded them. These 5 examples are further evidence that the ensembles of folding and unfolding pathways are one and the same, and that protein folding is a microscopically reversible process. However, in the other 2 cases, En-HD passed through remarkably similar unfolding and refolding transition states, but only a portion of the actual refolding pathway was similar to the unfolding pathway.

## Methods

### Molecular Dynamics Simulations

A total of 9 MD simulations are addressed in this paper at the following temperatures with simulation times in parentheses: 298K (100 ns), 323K (100 ns, 50 ns, 42 ns), 330K (100 ns, 100 ns), 373K (24 ns, 75 ns), 498K (20 ns, 60 ns), for a total of 0.67  $\mu$ s. All 4 of the 373K and 498K simulations have been described previously (9,12,13), as have the first 2 323K and the 298K simulations (14).

Both of the 330K simulations were performed using our in-house molecular dynamics package, *in lucem* Molecular Mechanics (*ilmm*) (15) with the Levitt *et al.* force field (16) using previously described protocols (17). The crystal structure (PDB ID: 1ENH) was minimized for 1000 steps and solvated in a box of F3C water molecules (18) such that there was at least 12 Å between the protein and the edge of the periodic box. The density was set to 0.985 g/mol in agreement with the experimentally-determined liquid-vapor coexistence curve for this temperature (19). 1000 steps of steepest descent minimization were performed on the water alone followed by 1 ps of dynamics. Next, the water and the protein were independently minimized for an additional 500 steps. Production simulations were performed for 100 ns allowing all atoms to move with structures written out every 1 ps. Long-range interactions were truncated after 8 Å using a force-shifted nonbonded cutoff. Our force-shifted cutoff method at this distance is the most effective treatment of long-range interactions based on computational savings, energy conservation, and ability to reproduce experimental results (20). The 323K/3 simulation followed the same protocol, except there was only 8 Å of padding between the protein and the edge of the box, the protein was minimized for 200 steps before adding water, and the simulation was run for 42 ns.

### C $\alpha$ RMSD Matrix and 3D MDS

All-vs.-all C $\alpha$  RMSD matrices were calculated to identify clusters of structures with similar conformations. Granularities were chosen to give 1000-5000 time points over the period of interest. The C $\alpha$  RMSD between each structure and every other structure was computed, resulting in a matrix with 1000<sup>2</sup>-5000<sup>2</sup> data points. Low C $\alpha$  RMSD boxes on the diagonal represent a period of time during which the protein stayed in a particular conformation. When these boxes lie off the diagonal, they indicate conformations of similar structure visited discontinuously in time. As described previously (13), the “core” (residues 8-53) was usually used to calculate C $\alpha$  RMSD, rather than the whole protein (residues 3-56). Since the fluctuations of the terminal residues are not indicative of the overall motion of the protein and introduce noise, the 5 residues at the N-terminus and 3 at the C-terminus were not included where specified.

Using the program R (21), multidimensional scaling (MDS) was performed to project the matrix down to 3 dimensions. This scaling results in a 3D plot in which each point represents a structure, and the distance between any two points is proportional to the C $\alpha$  RMSD between the respective structures. The points are connected in order of time for the period of interest. As with the matrix, a series of points close together indicates that the structures are similar. Using this plot, TSEs were chosen as the last structures leaving the extended native cluster for an unfolding event (5) or the first structures upon returning to the native cluster for a refolding event. The TSE was defined as the point of cluster exit and previous 5 ps for unfolding (5) and as the cluster reentry and succeeding 5 ps for refolding.

### HIII-core Distance Calculation

The distance between the closest backbone atoms in the C-terminus of HIII and the HI-HII scaffold in the crystal structure was chosen to represent the movement of HIII. The atoms chosen were the C $\alpha$  of Phe 20 and the backbone carbonyl C of Lys 52, and the distance was measured at 10 ps granularity. Since the 373K/1 simulation was so short, a granularity of 1 ps was used to give consistent sampling.

### Average Structures

Average structures were calculated using 100 ps granularity for the long N and N' time-spans. For TS structures, all 6 structures in the TSE were included in the average. The C $\alpha$  RMSD of the core residues (8-53) between the average structures was then calculated.

### HIII-core Contacts

A contact for a pair of residues was defined based on whether any one of the heavy atoms in the first residue was below a set cutoff of any heavy atom in the second residue. This cutoff was defined as 5.4 Å for carbon-carbon distances and 4.6 Å for all other atom pairs. For Figure 4, the calculation was taken over the time period of interest with 1 ps granularity, and the percentage of structures in which the two specified residues were in contact was reported for the average measurements. For the whole-simulation graphs (Figure 1c), each of the contacts is listed as a separate horizontal line, and if the contact was present at the given time point along the x-axis, a cross (+) was plotted. 10 ps granularity was used.

To choose which contacts to report, we identified residue pairs in which one member of the pair was in HIII and the other was not. Of these, the only pairs that were selected were those that were in contact at least 25% of the time in the native (298K) simulation.

### Calculation of S-values

The S-value is a semi-quantitative structure index that provides an overall measure of the secondary and tertiary structure for each residue of the protein (6). S-values agree well with experimental  $\Phi$ -values for a variety of proteins (6,9,22-25). The S-value is the product of the extent of native secondary structure ( $S_{2^\circ}$ ) and native and nonnative tertiary contacts ( $S_{3^\circ}$ ) present in a given structure relative to the number of contacts in the crystal structure. A value of 1 for S corresponds to native-like extent of structure in the TS, while a value of 0 suggests the residue is unstructured. As previously described (9),  $S_{3^\circ}$  was used in place of S for residues Phe 8, Leu 26, and Leu 40. For these three residues, side-chain interactions were maintained despite disorder in the main chain. Consequently, the product of  $S_{2^\circ}$  and  $S_{3^\circ}$  did not accurately represent the degree of structure retention.

### Protein- DNA Interactions

In order to generate a semi-quantitative measure of whether MD-generated En-HD conformers bind DNA, we measured distances between the DNA sugar-phosphate backbone and the HI-

HII helical hairpin using the crystal structure of En-HD bound to DNA (PDB ID: 3HDD). The crystal structure contains two nearly identical En-HD structures (core C $\alpha$  RMSD = 0.27 Å), so we selected the one bound to the ideal TAATTA sequence for all measurements (26). En-HD binds DNA primarily through residues in HIII (major groove) and the N-terminus (minor groove) (26,27). Since the N-terminus becomes structured only upon binding DNA (26), its conformation during our simulations should have no bearing on whether free En-HD is structured enough to bind DNA. Using Profit (28), the MD structures were first aligned to the DNA-bound structure based on a least-squares fit of the C $\alpha$  atoms in HIII, the DNA-binding helix. We selected a pairs of residues for the distance measurements, representative of one of two hydrogen bonds in the DNA-bound crystal structure that did not involve HIII or the N-terminus of En-HD. The atoms chosen for the measurement were the C $\alpha$  of Tyr 25 from En-HD and backbone P of Thymine 28 from the DNA. For measurements over time, 10 ps granularity was used. A period of time beginning 1 ns after the TS and continuing for 1 ns was selected to represent D for all 4 high-temperature unfolding simulations. Since there was not a full nanosecond of denatured time for most of the 4 lower-temperature simulations, the most denatured structure based on 3D MDS of the C $\alpha$  RMSD matrix was used.

## Protein Images and Figures

All protein images were rendered using VMD (29), C $\alpha$  RMSD matrices were made using R (21), and 3D MDS images were rendered with Chimera (30). Graphs were plotted and rendered in Gnuplot (31).

## Results

A total of 10 MD simulations were performed at 5 different temperatures (298K, 323K, 330K, 373K, and 498K). We describe the major conformational states of En-HD in each of the 10 simulations: N, N', TS, and D. When En-HD was in N, HIII was docked against the HI-HII scaffold. N' was characterized by a slight movement of HIII towards the N-terminus without losing many contacts or solvating the hydrophobic core. When HIII moved out and away from the HI-HII scaffold, exposing the hydrophobic core, En-HD was deemed to be in D. The protein was not necessarily unfolded in D, but it was not native nor, by definition, biologically active. Whenever the protein moved from N' to D or from D back to N', a TS was identified. These four states will be discussed further in the context of each simulation.

## Overview of Simulations

**298K**—En-HD remained folded in the native 298K simulation with a core C $\alpha$  RMSD of  $2.1 \pm 0.3$  Å (average  $\pm 1$  standard deviation) and an HIII-core distance (Phe20 C $\alpha$  – Lys52 carbonyl C, see Methods) of  $10.6 \pm 1.4$  Å for the first 80 ns of the 100 ns simulation (Figure 1a,b). During the final 20 ns of the simulation, the 9 C-terminal residues formed a  $\pi$ -helix, but HIII remained docked to the HI-HII scaffold. For this reason, the final 20 ns are not considered here.

**323K/1**—In this 100 ns simulation, En-HD stayed mostly in N but briefly moved to N' from 18-23 ns. It transiently moved to N' about 4 more times over the next 30 ns then stabilized in N for the remainder of the simulation (Figure 1). The protein did not populate D during this simulation, so there were no TSs identified.

**323K/2**—En-HD was in N for the first ~20 ns. At 19 ns, HIII moved ~10 Å towards the N-terminus, entering N' (Figure 3a). It remained in N' until 39 ns when there was a large jump in core C $\alpha$  RMSD and HIII-core distance, reflecting the undocking of HIII and entrance into D. HIII moved far enough away from the core (20 Å) for it to lose 10 of its 11 native core contacts (Figure 1c) and for the hydrophobic core to be solvated. This altered position was only transient, however, and HIII moved back to its position in N' a short 0.28 ns later. Over the next 1 ns,



HIII moved back to its N position where it stayed for ~3 ns. HIII then returned to its N' position transiently before the protein once again entered D at 43 ns. This transition was marked by another jump in core C $\alpha$  RMSD, HIII-core distance, and loss of contacts (Figure 1). After 0.16 ns, the protein returned to N', where it remained for the duration of the 50 ns simulation. The structures of En-HD in all four of its states are shown in Figure 2, and these structures are representative of those seen in the other 6 simulations.

**323K/3**—The first 5 ns of this simulation were spent mostly in N, after which En-HD stabilized in N' for another 5 ns. There was a transient movement to D at 10 ns, indicated by a spike in core C $\alpha$  RMSD and HIII-core distance. There was a transition back to N for 2 ns with a transient jump to N' during that time. Then at 13 ns, there was a more stable transition to N', interrupted only by a transient jump to D at 14 ns. En-HD remained in N' until the end of the simulation. Only the first 25 ns of the 42 ns simulation are considered here in order to focus on the transitions of interest.

**330K/1**—En-HD was stable in N for the majority of this simulation, with the exception of 43-53 ns (Figure 1). During this time, HIII unwound slightly at the N-terminus causing it to move towards the N-terminus by a register, much like the movement previously described for N'. Because the movement was due to unwinding rather than a simple loop movement, N' in this simulation had somewhat different properties (core C $\alpha$  RMSD HIII-core and contact pattern) than N' in the 323K simulations (Figure 1).

**330K/2**—This simulation began with 8 ns of N, moved to N' for 1 ns, returned to N for 8 ns, then shifted to N' again. At 27 ns, the protein moved from N' to D for 0.33 ns, then returned to N' at 28 ns. Again, at 44 ns, En-HD moved from N' to D and stayed there until 47 ns when it returned to N'. It remained in N' for the duration of the 100 ns simulation (Figure 1). Each of the three non-transient N' states here agree very well with N' in the 323K simulations as shown in Figure 2.

**373K/1**—En-HD moved from N to N' within 0.5 ns, then proceeded on to D. An unfolding TS was previously identified at 1.720 ns (9,12,13), and we identified a new refolding TS shortly after 5 ns. After En-HD returned to N', it remained there for the remainder of the 24 ns simulation. We will focus only on the first 10 ns of the simulation here. After refolding, N' did not agree very well with N' in the other simulations near the  $T_m$ , due to the fluctuating  $\pi$ -helical structure adopted by the C-terminus of HIII around 7 ns (data not shown). The temperature of this simulation was nearly 40K above the  $T_m$  of En-HD, so refolding is expected to occur only rarely and incompletely. We were lucky to observe such an event so that we can compare folding and unfolding within a continuous trajectory at high temperature.

### Native' State at Elevated Temperature

Before HIII moved far away enough from the HI-HII scaffold for En-HD to become “denatured,” it occupied two distinct positions. The conformation of En-HD with HIII positioned most similarly to the crystal structure is N, and we define N' as the state in which HIII slides ~10 Å towards the N-terminus (Figure 3a,b).

In all three of the 323K simulations and the second 330K simulation, the same N' structures were observed based on core C $\alpha$  RMSD and HIII-core contact similarity (Figures 3c, 4a, and 5a). Based on the core C $\alpha$  RMSD matrix of the first two 323K simulations (Figure 3c), the conformations from 18-23 ns and the subsequent transient deviations from N in 323K/1 were the same as those from 20-40 ns and 42-50 ns (with the exception of the unfolding events) in 323K/2. The core C $\alpha$  RMSD between average structures of N' in the 323K/1-3 and 330K/2 simulations also showed this similarity (Figure 5a). The 0.5-1.5 ns N' in 373K/1 was in good

agreement with N' in 323K/1-3 and 330K/2 based on core C $\alpha$  RMSD, but the 6-10 ns N' was less similar due to disruption of HIII.

The fraction of time the HIII-core contacts were made was remarkably similar for all 10 non-transient N' conformations, with the exception of 330K/1 (Figure 4a). The average standard deviation for contact time over all 11 residue pairs was 19%, and if the 330K/1 simulation was excluded, it dropped to 15% (Figure 4). Contact vs. time plots are useful to probe which pairs are in contact over the course of a simulation (Figure 1c). Based on this, contacts Phe 49 – Arg 24 (F), Phe 49 – Leu 26 (G), Lys 52 – Phe 20 (H), and Arg 53 – Arg 24 (K) were lost when En-HD moved from N to N' and were mostly regained if it moved back to N from N'.

### Properties of the Transition State Ensembles

A total of 13 new unfolding and refolding TSEs were identified in 4 simulations at 323K, 330K, and 373K. The TSEs all had low core C $\alpha$  RMSD to each other ( $2.1 \pm 1.0$  Å), particularly the  $T_m$  TSEs ( $1.7 \pm 0.9$  Å, Figure 5b). The  $T_m$  TSEs were less similar to the high-temperature TSEs, with lower core C $\alpha$  RMSDs to the 373K TSEs than the 498K TSEs, but most were 2-4 Å core C $\alpha$  RMSD between any two high temperature and  $T_m$  TSE average structures. The core C $\alpha$  RMSD to the native state was  $3.1 \pm 0.4$  Å over all 17 TSEs, and the lowest core C $\alpha$  RMSDs were observed between the exit and reentry TSEs at 39 and 43ns in the 323K/2 simulation (0.65 Å in both cases).

The HIII-core contacts agreed very well between the new and previously identified TSEs (Figure 4b). The Ile 45 – Leu 40 contact was consistently sustained in all of the TSEs, likely due to the residues' positions in the HII-HIII loop and at the N-terminal end of HIII, respectively. Where there were dissimilarities in contacts made between the known and new TSEs, it was usually the case that there were less contacts made in the new, lower temperature TSEs.

There was a pattern of gain and loss of HIII-core contacts preceding and following the TSEs across the simulations. There were 6 residues involved in contacts characteristically lost in N': Phe 20, Arg 24, Leu 26, Phe 49, Lys 52, and Arg 53; and aside from the 4 pairs that lost contact in N', these 6 residues were also involved in 3 more contacts: Phe 49 – Phe 20, Arg 53 – Phe 20, and Arg 53 – Leu 26. Of these 3 pairs, the 53-26 contact was lost early in all 6 simulations at elevated temperature (Figure 1c). In the case of the 12 new TSEs in the 323K and 330K simulations, the 49-20 contact was lost within 1 ns before or immediately after the exit TS, and it was reformed within 0.1 ns following the reentry TSs. The 53-20 contact was also lost during in the same time (with one exception: 323K/3 10 ns), and it was regained within 1 ns of all 4 reentry TSs in the 323K simulations, but it never reformed after the first reentry TS in the 330K/1 simulation. In the 373K/1 simulation the 49-20 pair was lost 0.5 ns after the exit TS and regained 0.1 ns after the reentry TS, and the 53-20 pair was lost 1.5 ns before the exit TS and regained 1 ns after reentry.

S-values, which quantify local structure in TSEs, were calculated for the 13 TSEs and compared with experimentally determined  $\Phi$ -values. The S- and  $\Phi$ -values agree very well for 10 of the 13 new TSEs, with linear correlation coefficients ranging from  $R = 0.71$ - $0.86$  (Figure 6). The two 330K/2 44-47 ns TSEs did not agree well with experiment (S vs.  $\Phi$ ,  $R = 0.10$  and  $0.03$ ). Loss of secondary structure in the termini of HIII explains some of the disagreement. For example, Ala 43 from the N-terminal end of HIII completely lost its secondary structure, giving it an S-value of 0 rather than its  $\Phi$ -value of 1.05. The 373K/1 5 ns reentry TS had a slightly lower correlation coefficient ( $R = 0.60$ ), with the largest disagreements seen for residues Leu 26, Leu 38, and Gly 39 which are located at the ends of HII. The secondary structure was as expected for these residues (turn, helix, helix; respectively), so the discrepancy is due to altered packing.



## Protein- DNA Interactions

When En-HD can no longer bind DNA, it loses its biological activity and is therefore denatured, although it may not be unfolded. Consequently, it is informative to determine whether the structures we identified as N, N', and D can bind DNA. In an effort to quantify En-HD's DNA binding ability, we fit En-HD to DNA based on the C $\alpha$  atoms of the major binding helix (HIII) and assessed to what extent the other binding interactions could be formed. Tyr 25 is not in HIII or the unstructured N-terminus, and it makes a hydrogen bond to the DNA backbone (26,27). Since we forced HIII to bind, this residue was selected for distance calculations. Across the simulations, the distances for Tyr 25 – Thymine 28 show a general trend of being longer for D than N, but the distance is also longer in N' than N. Structures from the 323K/2 simulation and the Tyr 25 – Thymine 28 distance are given in Figure 9 and are representative of what we saw in the other simulations.

## Unfolding and Refolding Pathways

There were 7 unfolding and refolding events identified in 4 simulations during the following times: 39 and 43 ns in 323K/2, 10 and 14 ns in 323K/3, 27-28 and 44-47 ns in 330K/2, and 1-5 ns in 373K/1. Comparison of the structures indicates that the protein moved through the same conformations when it refolded as when it unfolded for a given unfolding/refolding event. This path retracing can be seen as an "X" on the core C $\alpha$  RMSD matrix and as overlaid paths in the 3D projection of the matrix (Figure 7). As the protein moved from its most denatured point back to N', a line perpendicular to the diagonal of the matrix is apparent. This line represents a series of conformations that is very similar to the series of conformations the protein passed through just previously, but in reverse order. In the 3D projection, this same phenomenon is seen as overlapping points along the path from N' to the most denatured point, then back to N'. This evidence was present to different extents for each of the 7 unfolding/refolding events, but it is the most striking for 323K/2 39 and 43ns, 323K/3 10ns, 330K/2 27-28ns, and 373K 1-5ns.

While N' is not in the denatured ensemble, it is distinct from N, and thus there must exist a low-energy pathway to move between the two states. As with the unfolding and refolding pathways, there was an "X" on the core C $\alpha$  RMSD matrix when En-HD transiently moved from N' to N and back. In the 323K/3 simulation, there were 2 transient N $\rightarrow$ N' $\rightarrow$ N movements. There was an "X" visible on the core C $\alpha$  RMSD matrix around both N $\rightarrow$ N' $\rightarrow$ N transitions, and there was a third "X" off the diagonal, around the intersection of the times corresponding to both of the individual N $\rightarrow$ N' $\rightarrow$ N transitions (Figure 8). This third "X" indicates that not only were the N $\rightarrow$ N' and N' $\rightarrow$ N pathways very similar for a single transition, but also that the two N' $\rightarrow$ N $\rightarrow$ N' transitions were almost equally as similar.

## Discussion

Three different conformational states were populated by En-HD in our 7 simulations: N, N', and D. The protein's state was determined based on a combination of measurements: core C $\alpha$  RMSD, HIII-core distance, and HIII-core contact pattern (Figure 1). Measuring the core C $\alpha$  RMSD to the minimized crystal structure was the foremost method in determining En-HD's state. The core C $\alpha$  RMSD generally fluctuated between 1-3 Å when the protein was in N and 2-4 Å for N' (Figures 1a and 5). Values over 4.5 Å usually indicated a departure from N' to D. The Phe 20 C $\alpha$  – Lys 52 C distance, representative of the distance between HIII and the HI-HII scaffold, was also a good indicator of state. When the protein was in N, the distance fluctuated around  $10 \pm 2$  Å, while an increase to  $14 \pm 2$  Å indicated movement to N' (Figure 1b). Distances over 20 Å occurred when En-HD moved to D. Movement between states was more clearly discerned based on the HIII-core distance than on core C $\alpha$  RMSD. There was nearly always a clean jump in distance between the different states, which suggests that N and

N' are distinct states despite both being "native." The pattern of contacts between HIII and the rest of the protein also helped discriminate different states. The contact pairs selected for analysis were deliberately chosen to be good representatives of N. Jumps in core C $\alpha$  RMSD to 2-4 Å and HIII-core distance to ~14 Å, which were characteristic of N', coincided with loss of 4 contacts pairs: Phe 49 – Arg 24, Phe 49 – Leu 26, Lys 52 – Phe 20, and Arg 53 – Arg 24 (Figure 1c).

To estimate the likelihood that the 3 different En-HD conformations binds DNA, we fit structures from our MD simulations onto the DNA-bound crystal structure and took a distance measurement that might discern between native (N or N') and D. Representative structures from the 323K/2 simulation and distances are given in Figure 9. Even though both the N' and D structures were distinct from the crystal structure, N' conformations could more easily move back to N and be in a position to bind the DNA (and thus be biologically active) than those in D. For N'→N movement to occur, HI and HII would have to slide along the DNA and HIII so that they could dock against HIII as in N. This movement was what we saw for all N→N' transitions in our simulations without DNA. For a D structure to move to N, it would have to expel the water from its solvated hydrophobic core before HI and HII could dock back onto HIII and the DNA. Overall, D appears to be too distorted to function properly, and N' falls somewhere between N and D. While N' may be able to recover and clamp down on the DNA, it is also possible that it is too distorted and therefore inactive.

Having defined three different states for En-HD in 7 independent simulations, it is interesting to consider the variations within a state and how En-HD passes between them. N' was observed in all 6 of the elevated temperature simulations, but it was more similar in 5 of them (323K/1-3, 330K/2, and 373K/1) than it was in 330K/1 based on core C $\alpha$  RMSD and HIII-core contacts (Figures 4a and 5a). Also, the first period of N' in the 373K/1 simulation matched the first 4 simulations, while the second period was ambiguous. In all cases, HIII moved a register towards the N-terminus, away from the HI-HII turn, but N' in 330K/1 and the end of 373K/1 was somewhat different from N' in the other 4 simulations and in the beginning of 373K/1, despite having the same overall topology. This difference suggests that there may be multiple, subtly different N' states.

In the 4 simulations where there was N'↔D movement, a total of 13 new TSEs were identified. The TSEs within one temperature were most similar, and the best agreement was between the unfolding and refolding TSEs for a single transient unfolding/refolding event (Figure 5b). In these cases, only a small portion of D was sampled so it was likely that the protein would find a refolding path very similar to its unfolding path from the ensemble of paths available.

Based on the 11 HIII-core contact pairs selected for analysis, there was good agreement among all 17 TSEs (13 new and 4 previously identified, Figures 1c and 4b). In almost all of the cases where there was disagreement between the new and previously identified TSEs, it was the case that there were more contacts present in the high-temperature TSEs, which suggests they were more native-like. It is expected that high-temperature TSEs will more closely resemble N than TSEs at the protein's  $T_m$  due to Hammond effects. This phenomenon causes the structure of the TSE to become more native-like upon destabilization, in this case by heat (22, 32, 33).

Not only were the 13 new TSEs consistent with those previously identified, but 10 of them were in good agreement with experiment based on comparing calculated S-values to experimental  $\Phi$ -values ( $R = 0.71$ - $0.86$ , Figure 6). The correlation was somewhat lower ( $R = 0.60$ ) for the 373K/1 refolding TSE and significantly lower ( $R = 0.10$  and  $0.03$ ) for the 330K/2 44-47ns TSEs. The 330K/2 44-47ns unfolding/refolding event followed a cyclical path as it folded and unfolded, yet its unfolding and refolding TSEs were very similar to each other (average structure core C $\alpha$  RMSD = 1.77 Å, Figure 5b). The lack of agreement in this case

may be illustrating discrepancies between single-molecule behavior versus bulk measurements. That is, our aberrant TSE pair may be extreme members of the much broader ensemble probed experimentally.

The 13 new TSEs generally agreed well with each other, the 4 previously identified TSEs from high-temperature unfolding simulations, and experiment. The unfolding and refolding TSEs were equally similar, which is evidence that all 17 TSEs come from the same global ensemble of transition states for En-HD folding. Further, our data suggest that the ensemble of paths for unfolding and refolding is also the same, which supports our long-standing contention that protein folding is a microscopically reversible process.

The symmetrical order of contact pair loss and gain upon unfolding and refolding is further evidence that protein folding is microscopically reversible. The 6 residues involved in the 4 contact pairs that were characteristically lost in N' made 3 additional contacts (Figures 1c and 4). One pair was lost early in the simulations, and the other two, Phe 49 – Phe 20 and Arg 53 – Phe 20, were usually lost just before the unfolding TS and regained right after the refolding TS in the 7 unfolding/refolding events. These 6 residues make up the half of the HIII-core contacts closest to the C-terminus (Figure 4c). Arg 24, Leu 26, and Lys 52 lost all of their HIII-core contacts in N', but Phe 20, Phe 49, and Arg 53 maintained 2 of the original 7 contact pairs. It was not until right after these 2 contact pairs were lost that En-HD reached its TS and became denatured, and they reformed right after reentering N' from D in most cases. Thus, loss and gain of these 6 hydrophobic core contacts are critical steps on the  $N \leftrightarrow N'$  and  $N' \leftrightarrow D$  pathways, based on our 323 and 330K simulations. This is evidence for microscopic reversibility in protein folding because there is a consistent order of loss of contacts in unfolding that is repeated in reverse order upon refolding. However, the pattern of loss by these 6 contact pairs was not consistently repeated in the 4 high-temperature unfolding simulations previously run at 373 and 498K (data not shown).

The all-vs.-all core C $\alpha$  RMSD matrix and its 3D projection are arguably the best ways to observe the similarity of MD structures over time. Indeed, there was a visible “X” on the core C $\alpha$  RMSD matrix for 5 of the 7 unfolding/refolding events. Similarly, the structures from the unfolding TS to the most denatured point back to the refolding TS overlay on the 3D projection of the core C $\alpha$  RMSD matrix for these 5 unfolding/refolding events (Figure 7). Both the “X” and the overlaid paths indicate that the conformations En-HD moved through from the unfolding TS to the most denatured point had low core C $\alpha$  RMSD to the conformations En-HD took as it moved back to the refolding TS, but in reverse order.

“X”s were also visible for transient movements from N to N'. In 323K/3, not only was the  $N \rightarrow N'$  path similar to the  $N' \rightarrow N$  path, but both  $N \rightarrow N' \rightarrow N$  movements followed highly similar paths. There is an “X” on the diagonal of the all-vs.-all core C $\alpha$  RMSD matrix for each  $N \rightarrow N' \rightarrow N$  movement, and there is also an off-diagonal “X” at the intersection of the times corresponding to each of the  $N \rightarrow N' \rightarrow N$  events. While these four  $N \leftrightarrow N'$  pathways are not complete folding or unfolding pathways, they are transitions between discrete states along the full folding pathway.

En-HD never moved from N to D without first passing through N' each of the 7 times it unfolded. In fact, in the 323K/2 and 323K/3 simulations, it moved back to N between the two unfolding/refolding events, passing through N' on the way. Based on our simulations, the  $N \leftrightarrow N'$  paths were part of the same ensemble as were the  $N' \leftrightarrow D$  paths, and N' is a necessary step between N and D. Together, these findings are evidence that the entire  $N \rightarrow N' \rightarrow D$  and  $D \rightarrow N' \rightarrow N$  pathways are mirror images of the same process, and thus protein folding is a microscopically reversible process.

We identified and characterized four distinct states of En-HD: N, N', TS, and D which were consistent across simulations at 298, 323, 330, and 373K. Core C $\alpha$  RMSD, HIII-core contacts, HIII-core distance, and predicted DNA binding ability were used to discriminate between the states and place them on the folding pathway. We identified 7 transient denaturing events in 6 simulations and identified 13 new unfolding and refolding TSEs. The 13 new TSEs agreed well with 4 previously identified TSEs based on core C $\alpha$  RMSD and HIII-core contacts as well as with experimental data based on  $\Phi$ - and S-values. In 5 of the 7 transient denaturing events, the unfolding pathway was nearly identical to the refolding pathway. We also found two N $\leftrightarrow$ N' transitions that followed the same pathway in the folding and unfolding directions for both transitions. These phenomena are evidence that the ensemble of folding and unfolding pathways are one and the same and that protein folding can be a microscopically reversible process.

## Acknowledgments

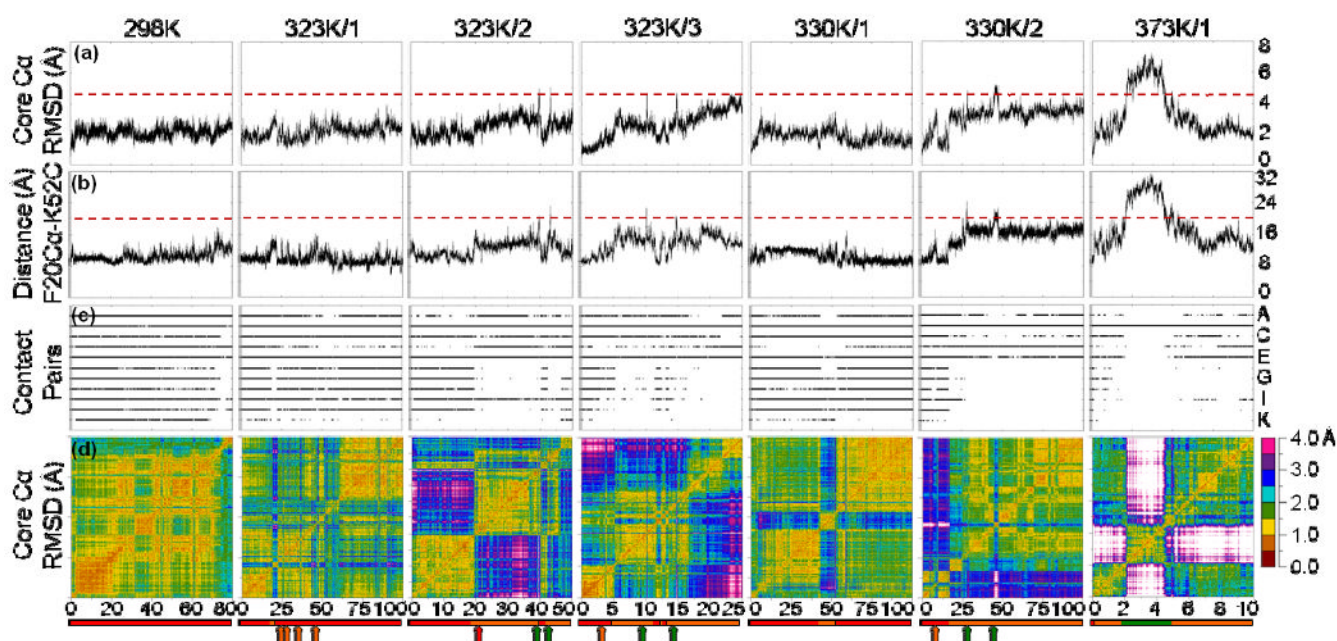
We thank Darwin Alonso, Amanda Jonsson, and Dustin Schaeffer for helpful discussions and technical assistance.

## References

1. Tolman RC. The Principle of Microscopic Reversibility. *Proc Natl Acad Sci U S A* 1925;11:436–439. [PubMed: 16587035]
2. Westheimer FH. Pseudo-rotation in the hydrolysis of phosphate esters. *Accounts of Chemical Research* 1968;1:70–78.
3. Minkin VI. Glossary of Terms Used in Theoretical Organic Chemistry. *Pure Appl Chem* 1999;71:1919–1981.
4. Jackson SE, elMasry N, Fersht AR. Structure of the hydrophobic core in the transition state for folding of chymotrypsin inhibitor 2: a critical test of the protein engineering method of analysis. *Biochemistry* 1993;32:11270–11278. [PubMed: 8218192]
5. Li A, Daggett V. Characterization of the transition state of protein unfolding by use of molecular dynamics: chymotrypsin inhibitor 2. *Proc Natl Acad Sci U S A* 1994;91:10430–10434. [PubMed: 7937969]
6. Daggett V, Li A, Itzhaki LS, Otzen DE, Fersht AR. Structure of the transition state for folding of a protein derived from experiment and simulation. *J Mol Biol* 1996;257:430–440. [PubMed: 8609634]
7. Day R, Daggett V. Direct observation of microscopic reversibility in single-molecule protein folding. *J Mol Biol* 2007;366:677–686. [PubMed: 17174331]
8. Itzhaki LS, Otzen DE, Fersht AR. The structure of the transition state for folding of chymotrypsin inhibitor 2 analysed by protein engineering methods: evidence for a nucleation-condensation mechanism for protein folding. *J Mol Biol* 1995;254:260–288. [PubMed: 7490748]
9. Gianni S, Guydosh NR, Khan F, Caldas TD, Mayor U, White GW, DeMarco ML, Daggett V, Fersht AR. Unifying features in protein-folding mechanisms. *Proc Natl Acad Sci U S A* 2003;100:13286–13291. [PubMed: 14595026]
10. Mayor U, Grossmann JG, Foster NW, Freund SM, Fersht AR. The denatured state of Engrailed Homeodomain under denaturing and native conditions. *J Mol Biol* 2003;333:977–991. [PubMed: 14583194]
11. Mayor U, Guydosh NR, Johnson CM, Grossmann JG, Sato S, Jas GS, Freund SM, Alonso DOV, Daggett V, Fersht AR. The complete folding pathway of a protein from nanoseconds to microseconds. *Nature* 2003;421:863–867. [PubMed: 12594518]
12. Mayor U, Johnson CM, Daggett V, Fersht AR. Protein folding and unfolding in microseconds to nanoseconds by experiment and simulation. *Proc Natl Acad Sci U S A* 2000;97:13518–13522. [PubMed: 11087839]
13. DeMarco ML, Alonso DOV, Daggett V. Diffusing and colliding: the atomic level folding/unfolding pathway of a small helical protein. *J Mol Biol* 2004;341:1109–1124. [PubMed: 15328620]
14. Beck DA, Daggett V. A One-Dimensional Reaction Coordinate for Identification of Transition States from Explicit Solvent Pfold-Like Calculations. *Biophys J* 2007;93:3382–3391. [PubMed: 17978165]

15. Beck, DAC.; Alonso, DOV.; Daggett, V. *ilmm*, in *luce*m Molecular Mechanics. University of Washington; Seattle: 2000-2008.
16. Levitt M, Hirshberg M, Sharon R, Daggett V. Potential-Energy Function and Parameters for Simulations of the Molecular-Dynamics of Proteins and Nucleic-Acids in Solution. *Computer Physics Communications* 1995;91:215–231.
17. Beck DAC, Daggett V. Methods for molecular dynamics simulations of protein folding/unfolding in solution. *Methods* 2004;34:112–120. [PubMed: 15283920]
18. Levitt M, Hirshberg M, Sharon R, Laidig KE, Daggett V. Calibration and testing of a water model for simulation of the molecular dynamics of proteins and nucleic acids in solution. *Journal of Physical Chemistry B* 1997;101:5051–5061.
19. Haar, L.; Gallagher, JS.; Kell, GS. *NBS/NRC Steam Tables: Thermodynamic and Transport Properties and Computer Programs for Vapor and Liquid States of Water in SI Units*. Hemisphere; Washington, DC: 1984.
20. Beck DA, Armen RS, Daggett V. Cutoff size need not strongly influence molecular dynamics results for solvated polypeptides. *Biochemistry* 2005;44:609–616. [PubMed: 15641786]
21. Team, R. C. R: A language and environment for statistical computing. R Foundation for Statistical Computing; Vienna, Austria: 2004.
22. Daggett V, Li A, Fersht AR. Combined molecular dynamics and phi-value analysis of structure-reactivity relationships in the transition state and unfolding pathway of barnase: Structural basis of Hammond and anti-Hammond effects. *J Am Chem Soc* 1998;120:12740–12754.
23. Li A, Daggett V. Molecular dynamics simulation of the unfolding of barnase: characterization of the major intermediate. *J Mol Biol* 1998;275:677–694. [PubMed: 9466940]
24. Fulton KF, Main ER, Daggett V, Jackson SE. Mapping the interactions present in the transition state for unfolding/folding of FKBP12. *J Mol Biol* 1999;291:445–461. [PubMed: 10438631]
25. Day R, Daggett V. Sensitivity of the folding/unfolding transition state ensemble of chymotrypsin inhibitor 2 to changes in temperature and solvent. *Protein Sci* 2005;14:1242–1252. [PubMed: 15840831]
26. Fraenkel E, Rould MA, Chambers KA, Pabo CO. Engrailed homeodomain-DNA complex at 2.2 Å resolution: a detailed view of the interface and comparison with other engrailed structures. *J Mol Biol* 1998;284:351–361. [PubMed: 9813123]
27. Kissinger CR, Liu BS, Martin-Blanco E, Kornberg TB, Pabo CO. Crystal structure of an engrailed homeodomain-DNA complex at 2.8 Å resolution: a framework for understanding homeodomain-DNA interactions. *Cell* 1990;63:579–590. [PubMed: 1977522]
28. Martin, ACR. *ProFit: Protein Least Squares Fitting*. University College London; London, England: 1992-2001.
29. Humphrey W, Dalke A, Schulten K. VMD: Visual Molecular Dynamics. *Journal of Molecular Graphics* 1996;14:33–38. [PubMed: 8744570]
30. Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, Ferrin TE. UCSF Chimera—A visualization system for exploratory research and analysis. *Journal of Computational Chemistry* 2004;25:1605–1612. [PubMed: 15264254]
31. Williams, T.; Kelley, C. Gnuplot. 1986-1993 1998 2004. [www.gnuplot.info](http://www.gnuplot.info)
32. Matthews JM, Fersht AR. Exploring the energy surface of protein folding by structure-reactivity relationships and engineered proteins: observation of Hammond behavior for the gross structure of the transition state and anti-Hammond behavior for structural elements for unfolding/folding of barnase. *Biochemistry* 1995;34:6805–6814. [PubMed: 7756312]
33. Day R, Bennion BJ, Ham S, Daggett V. Increasing temperature accelerates protein unfolding without changing the pathway of unfolding. *J Mol Biol* 2002;322:189–203. [PubMed: 12215424]

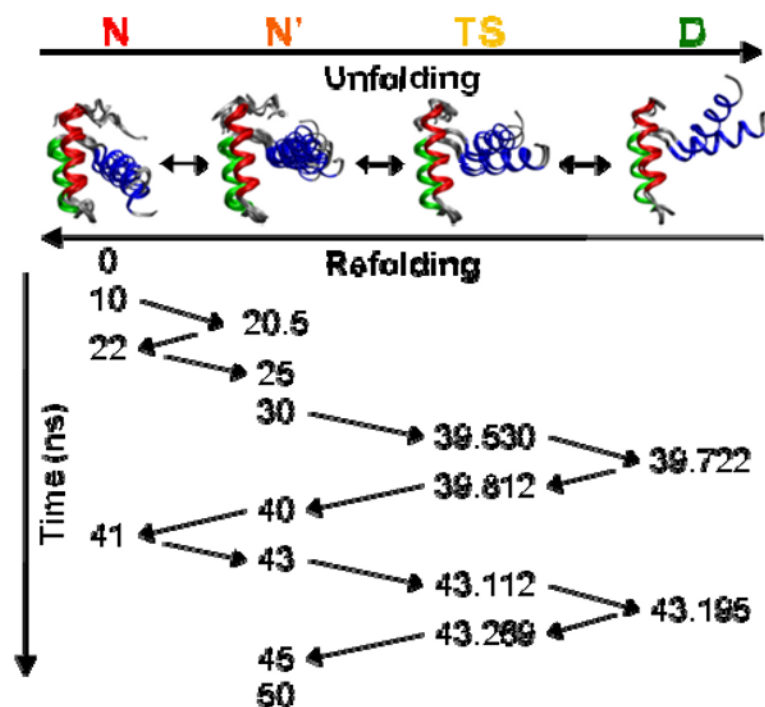




**Figure 1. General properties for each simulation**

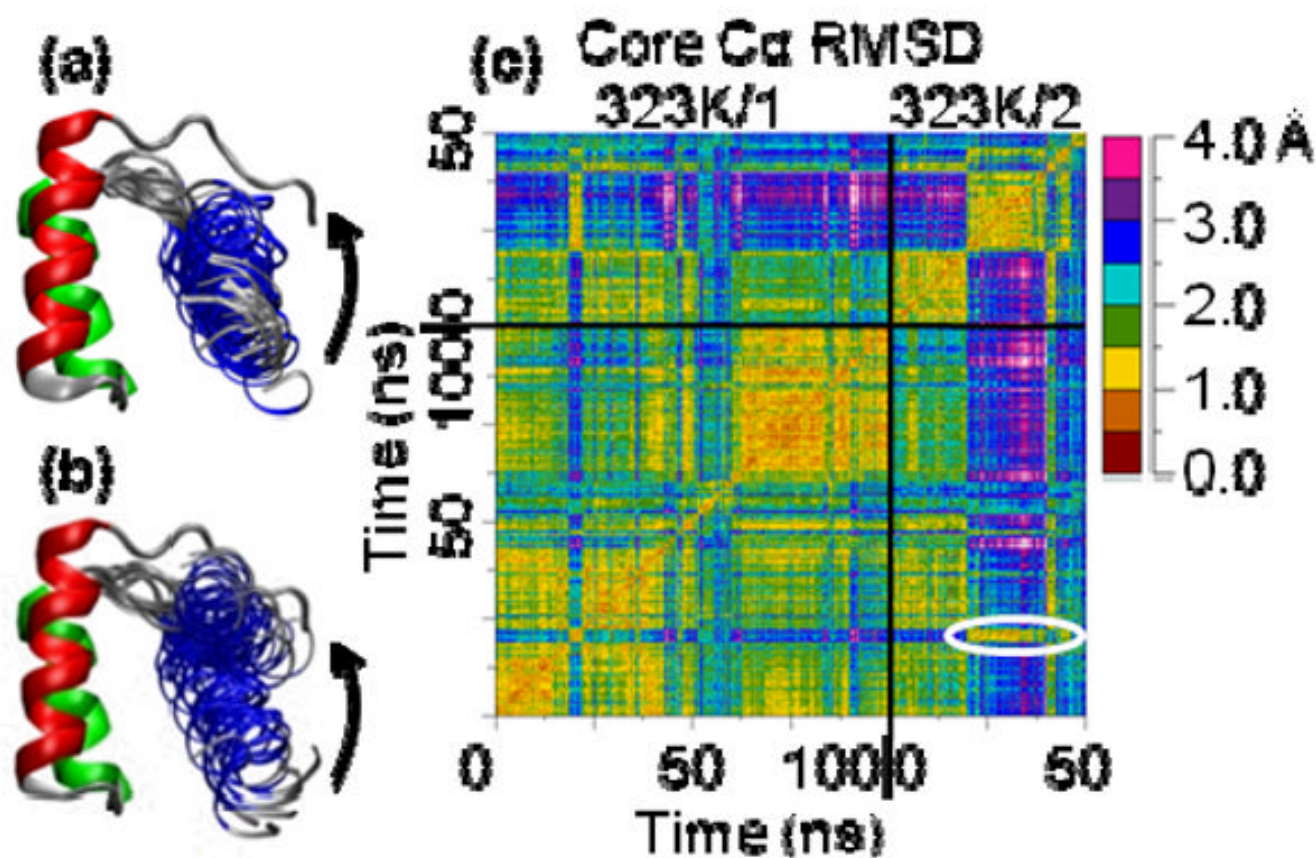
(a) C $\alpha$  RMSD of the core (residues 8-53) calculated over time for each simulation relative to the 0 ns structure. Values above  $\sim 4.5$  Å are indicative of movement to D, as indicated by the dashed red line. (b) Distance between the C $\alpha$  of Phe 20 and the backbone C of Lys 52 over time. N' is characterized by values of  $\sim 15$  Å, while distances of greater than 20 Å (dashed red line) appear when the protein moves to D. (c) Contacts made between residue pairs. Alternate pairs are labeled on the right from top to bottom: (A) Ile 45 – Leu 38, (B) Ile 45 – Leu 40, (C) Trp 48 – Leu 16, (D) Phe 49 – Leu 16, (E) Phe 49 – Phe 20, (F) Phe 49 – Arg 24, (G) Phe 49 – Leu 26, (H) Lys 52 – Phe 20, (I) Arg 53 – Phe 20, (J) Arg 53 – Arg 24, (K) Arg 53 – Leu 26. Contacts 49-24, 49-26, 52-20, and 53-24 are characteristically present in N but not N', while additional contacts are lost during D. (d) An all-vs.-all core C $\alpha$  RMSD matrix. Low- core C $\alpha$  RMSD squares on the diagonal represent a period of time with similar structures, and when they are off the diagonal, they indicate that the structures from the two corresponding time periods are similar. Below each matrix is a timeline depicting the different states the protein takes in each simulation: N (red), N' (orange), and D (green). Arrows represent transiently occupied states.





**Figure 2. Structures and order of population of the 4 states in the 323K/2 simulation**

Under each state is a set of representative structures from the course of the simulation overlaid and fit on HI-HII Ca atoms. HI (residues 8-20) is colored in red, HII (26-36) in green, and HIII (42-55) in blue. The time point in ns of each structure is listed below, and arrows connect them in the order they occurred. Structures within each state are similar, while each state is distinct. In N and N', HIII forms a  $\sim 15^\circ$  angle with the HI-HII scaffold. When En-HD reaches the TS, the angle is  $\sim 30^\circ$ , and it becomes even wider in D.

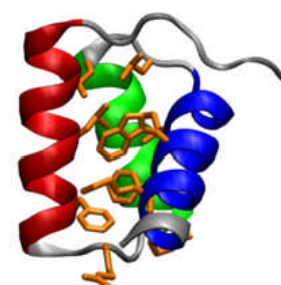


**Figure 3. N and N' structures and all-vs.-all core Cα RMSD matrix for 2 323K simulations**

(a) N and N' from the first 38 ns of the 323K/2 simulation taken every 0.5 ns for HIII. (b) N and N' from the first 100 ns of the 323K/1 simulation with structures taken every 1 ns for HIII. HI-HIII is from the 0 ns structure. HIII moves towards the N-terminus but not out, so the hydrophobic core is not solvated. (c) All-vs.-all core Cα RMSD matrix for the 323K/1 and 323K/2 simulations. 323K/1 and 323K/2 are in N' for the longest time from 19-20 ns and 20-39 ns, respectively. The color of the circled, low core Cα RMSD box off the diagonal indicates that these two N' states are as similar to each other as they are to themselves. Smaller boxes can be seen off the diagonal for the points in the 323K/1 simulation when En-HD moves to N' transiently indicating these transient N' states are the same as the longer two.

**(a) Fraction of Time in Contact: N and N'**

Res Pair	N	N'										N'	
	298K	323K				330K				373K			
		1	2	3		1	2		3		1	2	
		0-80	18-23	20-39	5-10	15-25	43-53	17-27	29-43	48-100	5-1.5	6-10	Avg
45-38	0.87	0.92	0.79	0.44	0.32	0.00	0.74	0.42	0.55	0.10	0.27	0.46	0.30
45-40	0.84	0.97	0.99	1.00	0.99	0.99	1.00	1.00	0.99	1.00	0.96	0.99	0.01
48-16	0.58	0.00	0.44	0.69	0.37	0.00	0.12	0.01	0.06	0.68	0.40	0.28	0.27
49-16	0.53	0.62	0.86	0.06	0.49	0.00	0.88	0.95	0.95	0.01	0.51	0.53	0.39
49-20	0.80	0.96	0.95	0.80	0.91	0.25	0.95	0.48	0.59	0.62	0.91	0.74	0.25
49-24	0.46	0.04	0.00	0.02	0.00	0.00	0.00	0.00	0.00	0.10	0.01	0.02	0.03
49-26	0.80	0.29	0.01	0.05	0.00	0.12	0.01	0.00	0.00	0.40	0.01	0.09	0.14
52-20	0.31	0.20	0.01	0.04	0.01	0.69	0.00	0.00	0.00	0.17	0.00	0.11	0.22
53-20	0.75	0.66	0.52	0.02	0.01	0.82	0.47	0.00	0.00	0.04	0.20	0.27	0.31
53-24	0.52	0.09	0.01	0.00	0.00	0.28	0.00	0.00	0.00	0.02	0.01	0.01	0.09
53-26	0.60	0.03	0.00	0.00	0.00	0.28	0.00	0.00	0.00	0.01	0.00	0.03	0.09
Avg	0.64	0.43	0.42	0.28	0.28	0.31	0.38	0.26	0.29	0.29	0.30	0.32	0.06
% 298	1.00	0.62	0.58	0.39	0.39	0.54	0.52	0.37	0.40	0.42	0.41	0.47	0.09

**(c)**

Leu 16    Ile 45  
 Phe 20    Trp 48  
 Arg 24    Phe 49  
 Leu 26    Lys 52  
 Leu 38    Arg 53  
 Leu 40

**(b) Fraction of Time in Contact: New and Known TSEs**

Res Pair	Known TSs Avg SD		323K/2				323K/3				330K/2				373K	New TSs Avg SD	
			39ns		43ns		10ns		14ns		27ns	28ns	44ns	47ns	5ns		
			U	R	U	R	U	R	U	R	U	R	U	R	R		
45-38	0.15	0.22	0.00	0.00	0.00	0.00	0.33	0.00	0.83	0.00	0.00	0.00	0.00	0.00	0.00	0.09	0.24
45-40	0.91	0.18	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.00
48-16	0.46	0.35	0.00	0.00	0.00	0.00	0.17	0.50	0.83	1.00	0.00	0.00	0.00	0.00	0.00	0.19	0.35
49-16	0.00	0.00	0.12	0.17	0.00	0.00	0.00	0.00	0.00	0.00	0.00	1.00	0.33	1.00	0.00	0.20	0.37
49-20	0.60	0.27	0.88	0.33	0.00	0.00	0.83	1.00	0.50	0.83	1.00	0.00	0.00	1.00	0.00	0.49	0.45
49-24	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
49-26	0.42	0.47	0.00	0.00	0.00	0.00	0.17	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.01	0.05
52-20	0.13	0.22	0.00	0.00	0.00	0.00	0.00	0.00	0.33	0.00	0.00	0.00	0.00	0.00	0.00	0.03	0.09
53-20	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
53-24	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
53-26	0.06	0.12	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Avg	0.25	0.17	0.18	0.14	0.09	0.09	0.23	0.23	0.32	0.26	0.18	0.18	0.12	0.27	0.09	0.18	0.07
% 298	0.35	0.26	0.23	0.17	0.11	0.11	0.28	0.30	0.48	0.36	0.22	0.28	0.16	0.39	0.11	0.25	0.12

**Figure 4. HIII-core residue pairs fraction of time in contact for N, N', and new TSEs in each simulation**

(a)-(b) are colored: < 25% pink, 25-75% purple, >75% blue. (a) Contacts for N in the 298K simulation and N' in the 323, 330, and 373K simulations. The temperature of the simulation, simulation number, and time span in ns of the N' state is given. (b) Contacts for new and previously identified TSEs. The simulation temperature and number is given (373K is 373K/1) as well as the ns during which the unfolding (U) or refolding (R) TSE occurred. The average fraction of time in contact was reported for each simulation as was the fraction of time in contact relative to the native (298K) simulation. Additionally, the average and standard deviation of the fraction of time in contact was calculated for each contact pair across all compared time spans. The average fraction of contacts is quite different for each of the states, N, N', and TS, with a low standard deviation, and those contacts that are lost are lost fairly consistently. (c) Each of the selected contact residues is shown on the En-HD structure in orange.

**(a) Core C $\alpha$  RMSD for N and N' Average Structures**

	N	N'									
		min	298K 0-80	1 18-23	2 20-39	3 5-10	1 15-25	2 43-53	3 17-27	2 28-43	1 48-100
N	min	0.00	1.82	2.22	2.72	2.41	3.10	2.20	2.68	3.10	3.22
	298K 0-80ns		0.00	1.69	2.48	2.23	2.82	2.26	2.41	2.62	2.70
	323K/1 18-23ns			0.00	0.96	1.22	1.60	1.77	1.01	1.57	1.58
	323K/2 20-39ns				0.00	1.57	1.27	2.26	0.60	1.39	1.28
	323K/3 5-10ns					0.00	1.68	1.74	1.66	1.72	1.86
	323K/3 15-25ns						0.00	2.54	1.36	1.35	1.26
	330K/1 43-53ns							0.00	2.38	2.61	2.67
	330K/2 17-27ns								0.00	1.30	1.28
	330K/2 28-43ns									0.00	0.40
	330K/2 48-100ns										0.00
N'	373K/1 0.5-1.5ns										0.00
	373K/1 6-10ns										0.00

**(b) Core C $\alpha$  RMSD for New and Known TSE Average Structures**

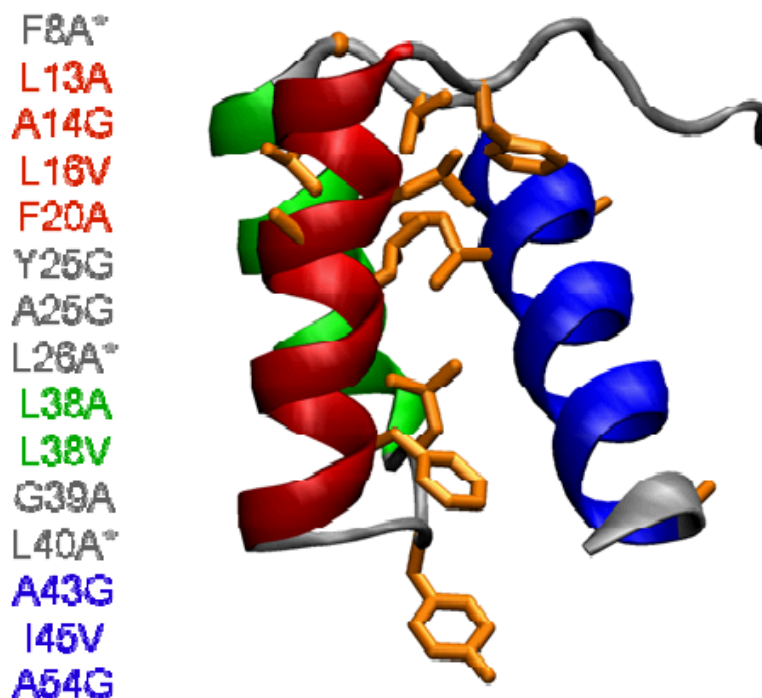
		N	New TS												Known TS				
		min	323K				330K				373K	373K		498K					
			39 U	39 R	43 U	43 R	10 U	10 R	14 U	14 R	27 U	28 R	44 U	47 R	1 5R	1 U	1 U	1 U	2 U
New TS	min	0.00	3.46	3.35	3.17	3.50	2.58	2.60	3.02	2.86	3.30	3.62	4.10	3.01	3.39	2.53	2.68	2.86	2.93
	323K/2 39ns Unfold		0.00	0.65	2.27	2.53	2.70	2.83	2.74	1.77	2.55	1.99	2.20	2.35	3.16	2.00	3.10	4.14	3.95
	323K/2 39ns Refold			0.00	2.09	2.32	2.49	2.63	2.45	1.41	2.42	1.86	2.21	2.24	3.00	1.89	3.02	3.84	3.63
	323K/2 43ns Unfold				0.00	0.65	1.31	1.35	1.74	1.84	1.11	1.74	2.14	2.01	1.91	1.70	2.22	2.72	2.90
	323K/2 43ns Refold					0.00	1.39	1.36	1.74	1.99	1.35	2.01	2.30	2.34	2.23	1.87	2.54	2.72	3.01
	323K/3 10ns Unfold						0.00	0.74	1.30	1.84	1.67	2.36	2.72	2.08	2.24	1.55	2.36	2.35	2.72
	323K/3 10ns Refold							0.00	1.56	1.92	1.81	2.51	2.79	2.25	2.41	1.46	2.14	2.09	2.63
	323K/3 14ns Unfold								0.00	1.57	2.18	2.41	2.71	2.23	2.70	1.92	2.50	2.65	2.97
	323K/3 14ns Refold									0.00	2.34	2.00	2.56	2.29	2.85	1.39	2.52	3.12	3.07
	330K/2 27ns Unfold										0.00	1.48	2.17	1.90	1.67	2.16	2.66	2.93	3.14
Known TS	330K/2 28ns Refold											0.00	1.75	1.96	2.03	2.37	3.00	3.66	3.60
	330K/2 44ns Unfold												0.00	1.77	2.65	2.60	3.14	3.93	3.95
	330K/2 47ns Refold													0.00	2.15	2.28	2.63	3.21	3.29
	373K/1 5ns Refold														0.00	2.71	2.94	3.08	3.06
	373K/1 1ns Unfold															0.00	2.24	2.84	2.92
	373K/2 1ns Unfold																0.00	2.45	2.80
	498K/1 0ns Unfold																	0.00	1.84
	498K/2 0ns Unfold																		0.00

**Figure 5. Matrix of core C $\alpha$  RMSDs between average structures representative of N, N', and TS**  
 The RMSDs reported are for the core C $\alpha$  residues of the average structure over the time periods indicated and are colored to visualize trends 0-1 Å (red), 1-2 Å (orange), 2-3 Å (yellow), 3-4 Å (green), 4-5 Å (blue). (a) Core C $\alpha$  RMSDs for average N' structures. The simulation temperature, number, and time span in ns of N' is given. Excluding the later 373K/1 N' structure, all of the N' structures are within 3 Å core C $\alpha$  RMSD of each other, and with the exception of 330K/1, all of the low-temperature N' structures are within 2 Å. (b) Core C $\alpha$  RMSD for the 13 new average TSE structures and 4 previously identified. The simulation temperature, number, and ns during which the unfolding (U) or refolding (R) TSE occurred is given. The average TSE structures are most similar to other TSE structures at the same temperatures, and all are about the same core C $\alpha$  RMSD from N (~2.5-3.5 Å).



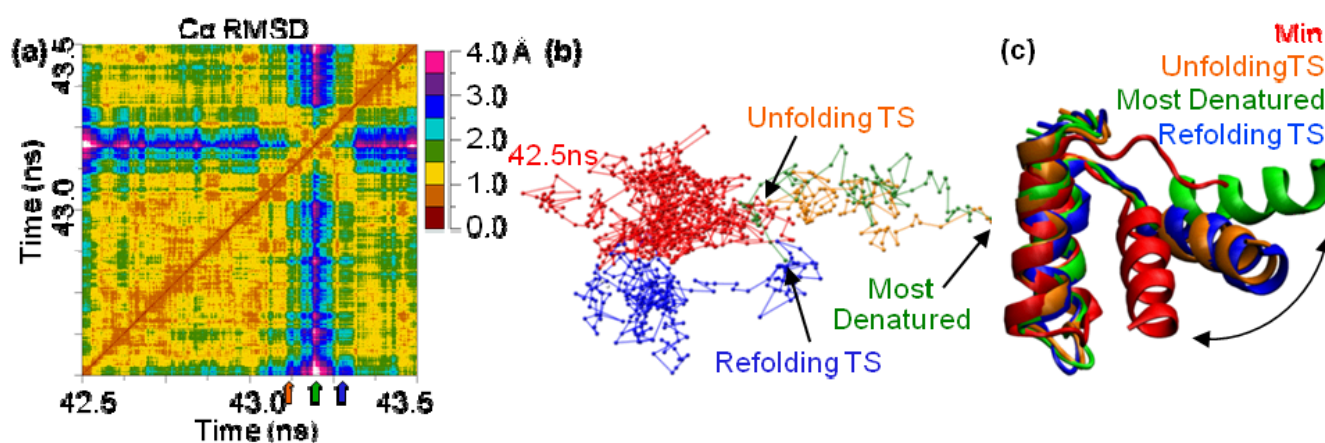
**Correlation Coefficients (R): S vs.  $\Phi$   
For Unfolding and Refolding TS Ensembles**

<b>323K/2</b>			
<u>39ns Unfold</u>	<u>39ns Refold</u>	<u>43ns Unfold</u>	<u>43ns Refold</u>
0.79	0.78	0.88	0.80
<b>323K/3</b>			
<u>10ns Unfold</u>	<u>10ns Refold</u>	<u>14ns Unfold</u>	<u>14ns Refold</u>
0.74	0.88	0.78	0.71
<b>330K/2</b>			
<u>27ns Unfold</u>	<u>28ns Refold</u>	<u>44ns Unfold</u>	<u>47ns Refold</u>
0.78	0.79	0.10	0.03
<b>373K/1</b>			
<u>5ns Refold</u>			
0.80			



**Figure 6. S-values for the 13 new TSEs**

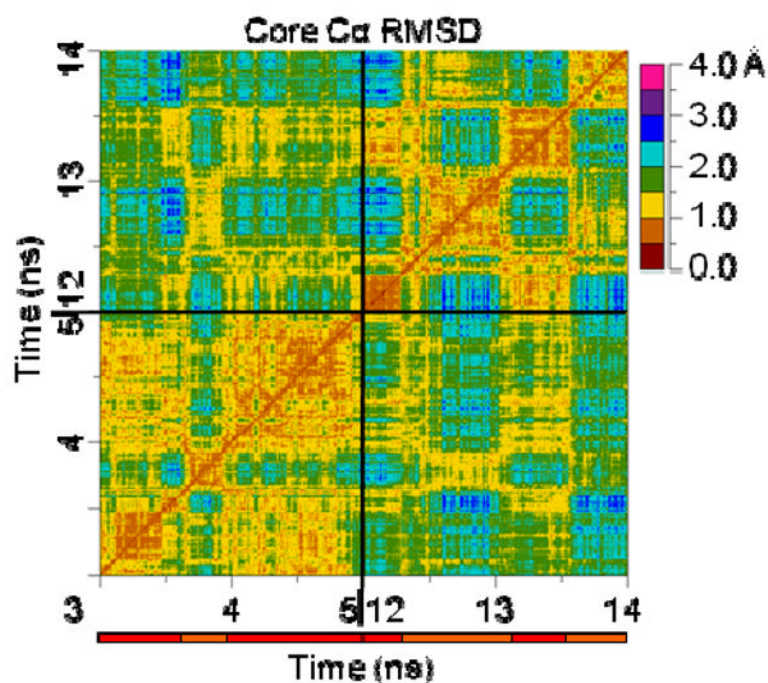
The correlation between calculated S-values and experimentally determined  $\Phi$ -values is quite good for the first 10 TSEs. The S- and  $\Phi$ -values do not agree as well for the 330K/2 44-47ns or 373K/1 TSEs, which is due to loss of secondary structure or altered packing. For the 3 residues noted by \*, only  $S_{3^\circ}$  was reported (see Methods).



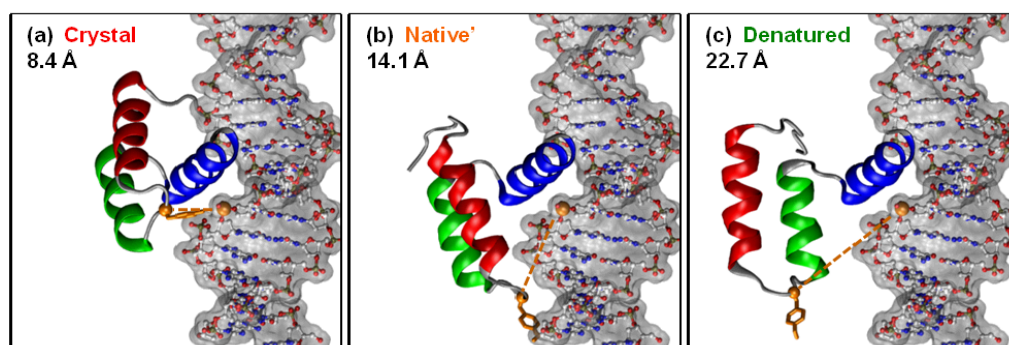
**Figure 7. Comparison of 323K/2 43 ns unfolding and refolding TSEs and pathway**

(a) An all-vs.-all Cα RMSD matrix showing the unfolding TS (43.112 ns, orange arrow), most denatured state (43.195 ns, green arrow), and the refolding TS (43.269 ns, blue arrow). The visible “X” is evidence that the conformations the protein takes as it leaves N’ to go to D are the same as those it takes when it returns to D, but in the reverse order. (b) The 3D MDS projection of the matrix from (a) in which each of the points represents a structure, and the distance between any two points is proportional to the Cα RMSD between the respective structures. The colors denote different periods in time: 42.5 ns to the unfolding TS (red), unfolding TS to the most denatured conformation (orange), most denatured to the refolding TS (green), and refolding to 43.5 ns (blue). That the paths that the protein followed as it moved from N’ to D and back are overlaid in the 3D projection indicates the conformations the protein took were very similar and in reverse order. (c) Structures of the unfolding TS (orange), refolding TS (blue), most denatured conformation (green), and the starting minimized structure (red). The structures were fit based on the Cα atoms of HI-HII. The 2 TS structures are nearly identical, while they are distinct from both N and D.





**Figure 8. All-vs.-all core C $\alpha$  RMSD matrix of 2 transient N $\rightarrow$ N' $\rightarrow$ N transitions from 323K/3**  
 Much like for the N' $\rightarrow$ D $\rightarrow$ N' transition in the 323K/2 43 ns TS, an “X” is visible on the core C $\alpha$  RMSD matrix for both N $\rightarrow$ N' $\rightarrow$ N transitions. In this case, there is a third “X” that is apparent off the diagonal at the intersection of the times at which the other two “X”s occur. This suggests that not only are the N $\leftrightarrow$ N' paths very similar, but also the N $\rightarrow$ N' $\rightarrow$ N path from the first transition is remarkably similar to that from the second.



**Figure 9. Structures from 323K/2 fit to DNA**

En-HD colored with Tyr 25 in orange sticks. En-HD's Tyr 25 C $\alpha$  and the DNA's Thymine 28 phosphate are shown in orange van der Waals spheres, and the distance between these two atoms is given and marked with an orange dashed line. (a) The crystal structure of En-HD bound to DNA (PDB ID: 3HDD) (b) N' at 30 ns in the 323K/2 simulation. (c) The most denatured state (D) from the 43 ns N'→D→N' transition, 43.195 ns. The structures show that En-HD moved slightly away from the DNA in N' compared to N, and much farther in D.