

1 Systematic comparison of Amber and Rosetta
2 energy functions for protein structure evaluation.

3 *Aliza B. Rubenstein^{1,2}, Kristin Blacklock^{3,4}, Hai Nguyen^{3,4}, David A. Case^{1,2,3,4*}, Sagar D. Khare^{2,3,4*}*

4 ¹Computational Biology & Molecular Biophysics Program

5 ²Institute for Quantitative Biomedicine

6 ³Department of Chemistry and Chemical Biology

7 ⁴Center for Integrative Proteomics Research

8 Rutgers, The State University of New Jersey, Piscataway NJ 08854

9

10 KEYWORDS: Amber, ff14SBonlySC, Rosetta, talaris2014, REF2015, comparative analysis,
11 structure prediction, benchmark, loop modeling

12

13

14

1 ABSTRACT An accurate energy function is an essential component of biomolecular structural
2 modeling and design. The comparison of differently derived energy functions enables analysis
3 of the strengths and weaknesses of each energy function, and provides independent benchmarks
4 for evaluating improvements within a given energy function. We compared the molecular
5 mechanics Amber empirical energy function to two versions of the Rosetta energy function
6 (talaris2014 and REF2015) in decoy discrimination and loop modeling tests. Both Rosetta's
7 talaris2014 and Amber's ff14SBonlySC energy functions performed well in scoring the native
8 state as the lowest energy conformation in many cases. In 24/150 cases with Rosetta, and in
9 2/150 cases using Amber, a false minimum is found that is absent in the alternative landscape. In
10 21/150 cases, both energy function-generated landscapes featured false minima. The newest
11 version of the Rosetta energy function, REF2015, which has more physically-derived terms than
12 talaris2014, performs significantly better, highlighting the improvements made to the Rosetta
13 scoring approach. To take advantage of the semi-orthogonal nature of these energy functions, we
14 developed a technique that combines Amber and Rosetta conformation rankings to predict the
15 most near-native model for a given protein. This algorithm improves upon predictions from
16 either energy function in isolation, and should aid in model selection for structure evaluation and
17 loop modeling tasks.

18 **Introduction**

19 Computational protein structure prediction is dependent on an accurate energy function. The
20 native state of a protein is expected to be found uniquely at the minimum of the energy function¹;
21 therefore, the energy function must robustly discriminate between native and non-native
22 conformations. A variety of energy functions to predict protein structure have been implemented
23 over the past forty years²⁻⁸. These potentials largely fall into one of two categories: molecular

1 mechanics force fields that rely on the combination of various empirical potentials such as
2 Lennard-Jones, torsional energies, Coulombic interactions, and desolvation penalties^{3,4,7} and
3 statistical or knowledge-based potentials that depend on characteristics of known protein
4 structures^{2,5,6}. While molecular mechanics force-fields are generally parameterized on small
5 molecule properties^{7,9-11}, statistical potential parameter optimization is often guided by known
6 biomolecular structures¹²⁻¹⁴. Each approach has its own drawback: since parameters in physically
7 derived force-fields are fit based on small molecule properties, they may not be suited to
8 macromolecules^{15,16}: for example, force-fields will often display biases towards secondary
9 structure propensities^{15,17}. On the other hand, statistical potentials are trained on specific datasets
10 of large biomolecules, and data sparseness may lead to overfitting¹⁸.

11 The Rosetta macromolecular modeling program energy function combines elements of both
12 categories; it contains physical force-field terms (Lennard-Jones interactions, electrostatic
13 interactions, desolvation penalties, etc.) and statistical potentials (probability of amino acid
14 identity given backbone angles, probability of backbone angles given amino acid identity,
15 probability of backbone-dependent rotamer, etc.)⁹. The most recent Rosetta energy function
16 (REF2015) is parameterized on both small molecule properties and large sets of biomolecular
17 structures¹⁸, although previous energy functions were generally parameterized on known
18 biomolecular structures alone¹³. While efforts have been made to compare the performance of
19 various empirical force-fields^{17,20,21}, little attention has been focused on the comparison between the
20 Rosetta energy function and empirical force-fields.

21 The Amber ff14SOnlySC force field¹⁰ uses a standard fixed-charge molecular mechanics
22 potential, with torsion potentials based entirely on fits to quantum chemistry data. It is very like
23 the more commonly-used ff14SB protein force field, but does not include the empirical

1 modifications to backbone torsion potentials that are present in ff14SB, and which provide an
2 improved balance of secondary structure in explicit solvent simulations. Hence, ff14SBOonlySC is
3 more "physics-based" than is ff14SB, and it arguably better suited for the implicit solvent
4 simulations used here, since the empirical backbone torsional potentials in ff14SB might be
5 specific to its use of explicit solvent simulations. The ff14SBOonlySC force field, in combination
6 with a generalized Born implicit solvent model², has been shown to fold a variety of single-
7 domain proteins using unrestrained molecular dynamics simulations³.

8 Comparing the Amber force-field and Rosetta energy function performance at structure
9 evaluation elucidates the strengths and areas of improvements for each energy function. As
10 Rosetta energy functions have been developed based on improving performance for certain
11 modeling datasets, testing their performance on the same macromolecular datasets may result in
12 overfitting of the Rosetta energy function, while comparing their performance to that of a
13 physics-based Amber energy function is a relatively unbiased comparison for evaluating
14 performance improvements. Finally, selecting a correct near-native model for a given sequence
15 is an elementary challenge; the combination of these two semi-orthogonal energy functions
16 provides a method for model selection that is able to select more accurate models.

17 **Methods**

18 **Benchmark Sets**

19 To evaluate and compare the performance of Rosetta and Amber energy functions, we used
20 two benchmark sets, a structure evaluation (decoy discrimination) set and a loop modeling set.
21 The decoy discrimination benchmark set includes a total of 150 proteins, a combination of two
22 independent decoy sets used in previous studies^{18,24}. The proteins in the set are monomeric and
23 have crystallographic native structures available in the RCSB PDB²⁵ with resolution $< 2.0 \text{ \AA}$. The

1 protein lengths range from 50 to 200 residues and have a diverse range of topologies. The decoy
2 sets were originally generated using biased and unbiased ab-initio sampling runs²⁶ followed by
3 parallel loophash sampling (PLS)²⁷. This produced 40,000-200,000 decoys per protein, ~1000
4 representative low-energy structures of which were chosen for each protein to cover the range of
5 possible C- α RMSD values.

6 The loop modeling benchmark set consisted of the 45-PDB dataset for 12-residue loops in the
7 monomeric protein loops training set of the 2016 Collaborative Assessment and Development of
8 Rosetta Energetics and Sampling (CADRES). This loop modeling benchmark set was obtained
9 from Shane O'Connor and Tanja Kortemme (personal communication).

10 **Structure Preparation**

11 **Rosetta**

12 In the decoy discrimination benchmark, the native crystallographic structures for each protein
13 set were downloaded from the RSCB PDB and residues were trimmed from the structure to
14 match the sequence of the crystal with the decoy structure in the benchmark sets. Native
15 structures were necessary to evaluate RMSD from native for decoy conformations. Native
16 structures were then relaxed using FastRelax²⁶ with the talaris2014¹³ scorefunction to relieve any
17 clashes. One hundred relaxation trajectories were simulated to generate one hundred relaxed
18 native-like decoys. These native-like decoys were used for false minima analysis. Then, these
19 one hundred native-like decoys, along with the ~1000 pre-sampled decoys, were subjected to
20 backbone and sidechain minimization using talaris2014 and the Limited-memory Broyden-
21 Fletcher-Goldfarb-Shanno (LBFGS) minimizer implementation with inexact line search
22 conditions (lbfgs_armijo_nonmonotone) over a maximum of 2000 iterations for convergence. C-
23 α atom RMSD was calculated for all decoys.

1 The REF2015 dataset was obtained from F. DiMaio and H. Park¹⁸. For this dataset, each decoy
2 was relaxed with 3 cycles torsion-space minimization and 2 cycles Cartesian mode²⁴ using the
3 REF2015 energy function⁹. Only 140 out of 150 protein systems were included in this set due to
4 the lower quality of experimentally determined structures for 10 systems (H. Park and F.
5 DiMaio, personal communication, July 5, 2017). Those 10 systems are ignored when comparing
6 REF2015 to Amber.

7 In the loop modeling benchmark, the native crystal structures for each protein set were
8 downloaded from the RCSB PDB and trimmed of excess residues that were not found in the
9 decoy PDB structures. The backbone and sidechain geometries for residues in the loop region of
10 each decoy structure were minimized in Rosetta using the talaris2014 scorefunction and the
11 lbfgs_armijo_nonmonotone over a maximum of 2000 iterations for convergence. C- α RMSDs
12 were calculated with respect to the crystal structure over loop residues only without fitting; since
13 the protein scaffold was fixed during optimization, this statistic describes the extent of loop
14 deviation. Loop residues are defined in supplementary file LoopDefs.xlsx.

15 **Amber**

16 Hydrogens were removed from the crystal structures and decoy PDBs, and initial structures
17 were built using the tLEaP module of AmberTools²⁸ with the ff14SBoonlySC¹⁰ forcefield
18 parameters. Minimizations were carried out for a maximum of 1000 steps under the LBFGS
19 quasi-Newton algorithm²⁹ with a convergence criterion of 0.01 kcal/mol-A. In the loop modeling
20 benchmark, positional restraints were added to all non-loop-residue atoms except for hydrogens
21 with a force constant of 10.0 kcal mol⁻¹ A⁻². Solvent effects were treated with a generalized
22 Born implicit solvent model (GB-Neck2²²) implemented in the Amber16²⁸ package with mbondi3
23 radii and a cutoff value of 999A for nonbonded interactions. Total potential energies of

1 minimized structures and C- α RMSDs with respect to the crystal structure were obtained using
2 the pytraj 2.0.0 interactive molecular dynamics simulation data analysis Python package³⁰, which
3 is a Python interface for cpptraj in AmberTools16³⁵. In the loop modeling benchmark, C- α
4 RMSDs were calculated over the loop residues only. Six sets of decoy structures for the loop
5 modeling benchmark were unable to be minimized in Amber due to missing residues, and those
6 sets were not considered in subsequent analyses (1cb0, 1dts, 1m3s, 1ms9, 1t1d, and 2pia).

7 **Energy Landscape Generation**

8 Energy landscapes (RMSD vs. normalized energy scatterplots) were generated for all proteins
9 for both Rosetta and Amber. The ideal shape of an energy landscape is that of a funnel (i.e.
10 Figure 1A, turquoise plot) where the lowest-scoring decoy conformations are of near-native
11 RMSD. We use the binned Boltzmann metric (see below) to evaluate the funnel shape of each
12 energy landscape.

13 **Energy Normalization**

14 For each set of energies per scorefunction per protein, energies are normalized so that the gap
15 between the 5th percentile of and the 95th percentile is equal to 1. This enables the comparison of
16 energies between different structures and between different energy functions. This is
17 accomplished via the following equation:

$$18 \quad E_{i(norm)} = (E_i - E_{min}) / (E_{95th} - E_{5th})$$

19 E_i refers to the raw energy of decoy i . E_{min} is the minimum energy value, E_{95th} is the 95th
20 percentile energy, and E_{5th} is the 5th percentile energy

21 **Funnel Evaluation Metric**

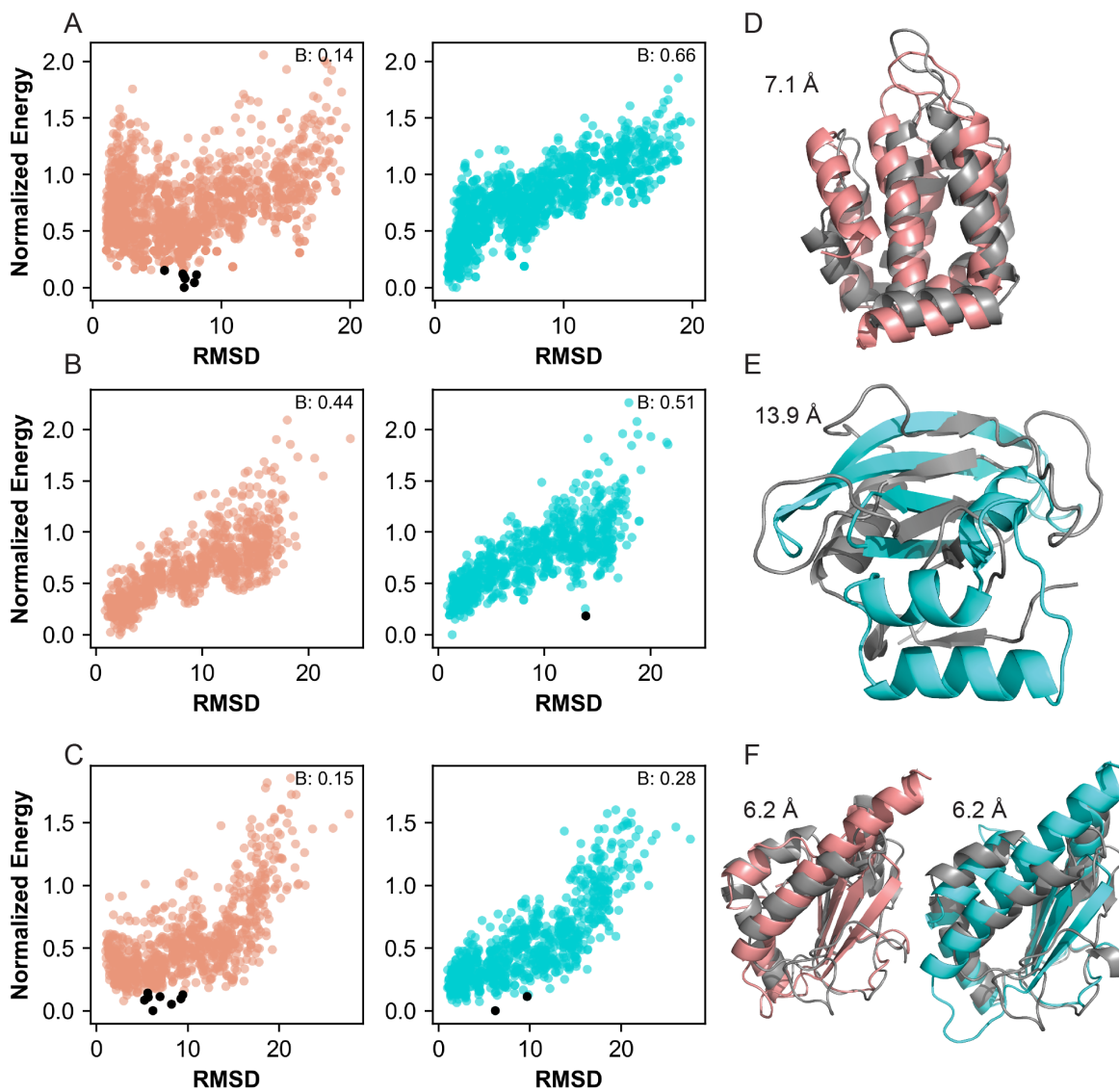
22 We use the binned Boltzmann metric, B , for energy landscape evaluation, as described
23 previously by Park et al.¹⁸. This metric finds the Boltzmann probability of selecting native-like

1 decoys over high-RMSD decoys based on their energy values. As in previous work²⁴, the metric
2 is averaged over multiple thresholds for determining native-like status for each decoy.

3
$$B = \frac{\sum_j (\sum_i d_{ij} P_i / \sum_i P_i)}{N_j}$$

4

5
$$P_i = e^{-\beta E_i(\text{norm})}$$



6

1 **Figure 1.** (A-C) Energy landscapes for 2QY7, 1T2I, and 1SEN respectively. Each dot on the
2 plot represents one decoy conformation. The x-axis is RMSD from native and the y-axis is
3 normalized energy. False minima (defined as decoys within top 10 energies but with RMSD >
4 5.0 Å) are depicted in black. The B metric, which represents the efficacy of the score-function at
5 differentiating between native and non-native decoys, is shown at the top right corner of each
6 plot. Rosetta plots are to the left, in salmon, and Amber plots are to the right, in turquoise. (D-
7 F) Superimposed native (gray) and Rosetta lowest-ranking false minimum decoy (salmon) and/or
8 superimposed native (gray) and Amber lowest-ranking false minimum decoy (turquoise) for
9 2QY7, 1T2I, and 1SEN respectively.

10 The conformation index is i and j is the native threshold definition index. Cutoffs are 0.5, 1.0,
11 1.25, 1.5, 1.75, 2.0, 2.25, 2.5, 2.75, 3.0, 3.5, 4.0, 5.0, 6.0 Å and N_j is thus 14. $E_{i(norm)}$ is the score of
12 decoy i as determined in Rosetta or Amber and normalized as described above. The factor β
13 refers to the Boltzmann factor and has a value of 0.1, as this was the value of β used by Park et
14 al. previously¹⁸. d_{ij} determines whether decoy i is considered native at threshold j ; it is set to 1 if it
15 is native and 0 if it is not. As the sum of the probabilities of the non-native-like conformations
16 approaches 0, the numerator ($\sum_i d_{ij} P_i$) approaches the value of the denominator ($\sum_i P_i$), so that
17 the value of B approaches to 1. As mentioned in Park et al.¹⁸, the B metric is better than the
18 previously used S metric²⁴ due to a larger increase in the metric for a poor energy landscape vs. a
19 good energy landscape than the increase from an already good energy landscape to a steeper
20 energy landscape. Additionally, it is a smoother metric that is less affected by single-decoy
21 outliers.

22 **Model Selection**

1 For each protein, a model was selected by finding the decoy that had the lowest sum of Amber
2 and Rosetta ranks; this decoy also satisfies the criteria for Pareto-optimality. First, the Amber
3 scores and Rosetta scores were converted into ranks so that the rank of decoy *a* was less than the
4 rank of decoy *b* if the energy of decoy *a* was less than the energy of decoy *b*. Second, the Pareto-
5 optimal solutions are found as follows. Decoy *a* is defined as dominating decoy *b* if both ranks
6 (Rosetta and Amber) of decoy *a* are \leq both ranks of decoy *b*. Pareto-optimal decoys are decoys
7 that dominate at least one other decoy and are not dominated by any decoys. From among the
8 set of Pareto-optimal decoys, the decoy that has the lowest sum of ranks is chosen as the
9 solution. In the rare case that more than one decoy has a minimum sum of ranks, a decoy is
10 arbitrarily chosen from the minimum-sum-ranks decoys. Technically, the minimum-sum decoy
11 must be found within the Pareto set of solutions; however, the use of the Pareto-minimization
12 allows for easier visualization and interpretation of the minimum-sum solution. Additionally, the
13 Pareto-optimal set of decoys may be useful for selecting a set of top-scoring *n* ($n > 1$) decoys
14 when performing consensus scoring according to the two energy functions.

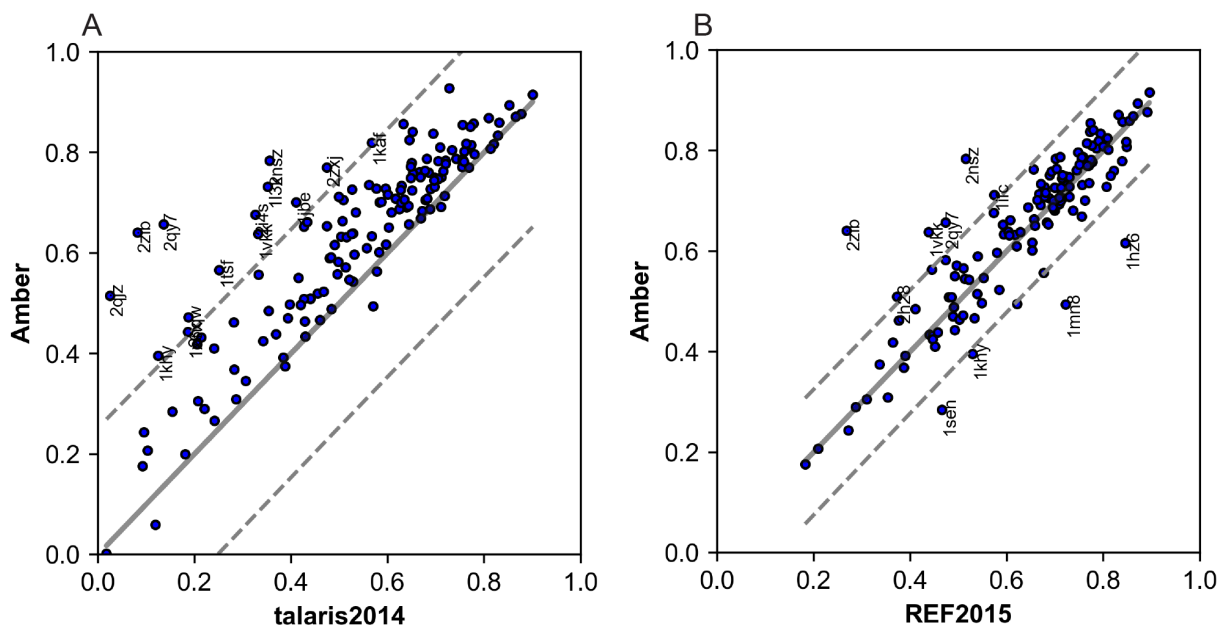
15 **Results**

16 **Performance of Amber and Rosetta energy functions in discriminating between native** 17 **and non-native structures**

18 Protein free energy landscapes involve folding funnels³¹⁻³³ which enable the folding chain to
19 efficiently find the native state³¹, and their existence implies that the higher energy of non-native
20 (decoy) structures compared to the native (e.g., crystallographically-determined) structure drives
21 protein folding. Therefore, a common test used for evaluating^{23,26} and improving¹⁸ energy
22 functions is the decoy discrimination test, in which the evaluated scores of decoy structures are
23 compared to that of near-native structures. High-RMSD decoys which have comparable energies

1 to near-native structures are classified as “false minima”, and are indicative of inaccuracies in the
2 energy function. The B metric¹⁸, ranging from 0 to 1, quantifies the existence of false minima in
3 a set of structures upon evaluation with a given energy function, with values close to 1 indicating
4 a smooth folding funnel with no false minima. Conversely, a lower B value indicates that one or
5 more false minima exist.

6 We compared the performance of the Amber energy function and Rosetta energy function at
7 ranking native state structures lower than decoy conformations for a set of 150 proteins. Amber
8 ff14SBoonlySC generally performed better than Rosetta talaris2014, scoring significantly higher
9 B metrics for many systems (Figure 2A). We also compared Amber to the newer default Rosetta
10 energy function, REF2015⁹, and found that while Amber did have a higher B metric for several
11 systems, several other systems had a higher B metric when scored by REF2015, thus showing the
12 improvement of REF2015 over talaris2014 when compared to Amber as an unbiased benchmark.
13 Nonetheless, the comparative performance of the two energy functions (Amber and REF2015;
14 Fig. 2B) shows that each has its strengths and limitations (Table 1). Our analysis was carried out
15 with the talaris2014 energy function, and we refer to it as the Rosetta energy function in the
16 remainder of this paper.



1
 2 **Figure 2.** Scatterplots to depict general performance of Rosetta talaris2014 scoring function vs.
 3 Amber scoring function (A) and Rosetta REF2015 scoring function vs. Amber scoring function
 4 (B) over entire decoy discrimination set. Each dot represents the B metric for one system. The
 5 black line is $x=y$ and the dashed line represents the 95% prediction interval. Any points that lie
 6 outside the 95% prediction interval are annotated with the PDB ID of that system.

7 We examined cases in which either Amber, Rosetta, or both were unable to correctly rank
 8 high-RMSD decoy conformations, scoring them as low-scoring instead of high-scoring. A false
 9 minimum is defined as a decoy within the top-10 ranked decoys that has a C- α RMSD from
 10 native of greater than 5 Å. Three of these cases are shown in Figure 1. 2QY7 (Figure 1A,D) has
 11 several false minima for Rosetta but none for Amber. Generally, Rosetta alone had at least one
 12 false minimum in 16% of structures. 1T2I (Figure 1B,E) has a false minimum for Amber but
 13 none for Rosetta; 1.32% of systems have at least one false minimum for Amber alone. 1SEN

1 (Figure 1C,F) has false minima for both Amber and Rosetta, as do 14% of overall structures
 2 (Table 1).
 3 **Table 1.** *B* metric, false minima, and model selection summary comparisons for Amber
 4 ff14SBonlySC, Rosetta talaris2014, and Rosetta REF2015 energy functions.

	No. Cases/Total No. Proteins
Decoy Discrimination	
ff14SBonlySC B > talaris2014 B by 0.1	54/150
talaris2014 B > ff14SBonlySC B by 0.1	0/150
ff14SBonlySC B > REF2015 B by 0.1	6/140
REF2015 B > ff14SBonlySC B by 0.1	9/140
False minima in ff14SBonlySC only (not talaris2014)	2/150
False minima in talaris2014 only (not ff14SBonlySC)	24/150
False minima in ff14SBonlySC and talaris2014	21/150
False minima in REF2015 only (not ff14SBonlySC)	0/140
False minima in ff14SBonlySC and REF2015	10/140
Minimum-sum RMSD < ff14SBonlySC selected RMSD by 1 Å	10/150
Minimum-sum RMSD < talaris2014 selected RMSD by 1 Å	21/150
Minimum-sum RMSD < ff14SBonlySC selected and talaris2014 RMSD by 1 Å	1/150
Loop Modeling	
ff14SBonlySC B > talaris2014 B	15/39
talaris2014 B > ff14SBonlySC B	7/39

5 Superimpositions of false minima decoys with native decoys show their distinct non-native
 6 conformations involving both misprediction of secondary structure elements as well as their

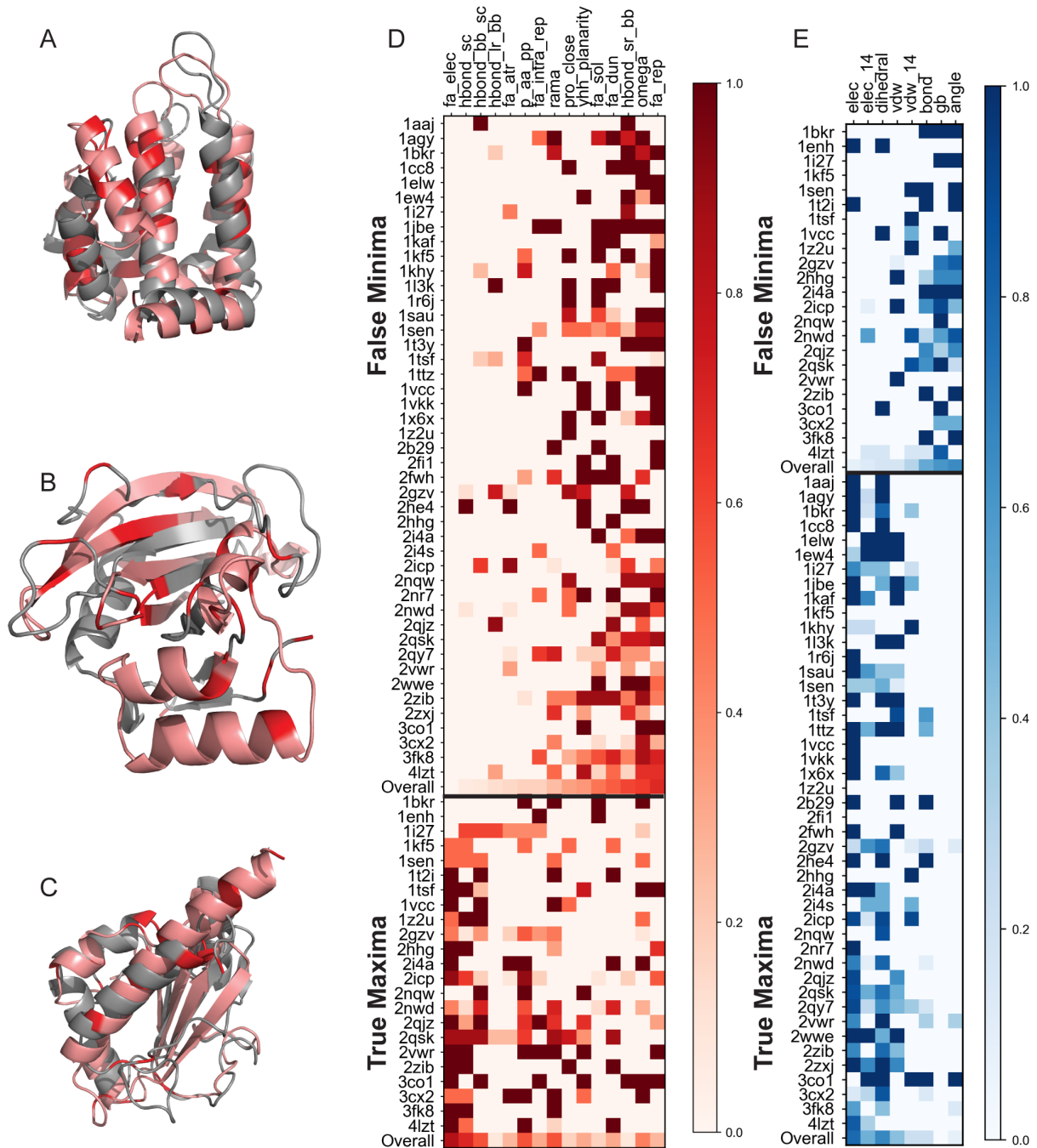
1 incorrect relative placement in tertiary structures. In the case of 2QY7, a Rosetta false minimum,
2 the four helical bundle found in the native structure is perturbed in the false minimum, as the
3 order of the first two helices is reversed; thus, they do not contact the other two helices as tightly
4 as that of the native structure (Supplementary Figure 1E, I-J). The difference between the native
5 structure of 1T2I and its Amber false minimum is more subtle. While the contact maps for the
6 native and false minimum conformations are similar, except for a small contact region in the
7 native structure between residues 40 and 59 that does not appear in the false minimum
8 (Supplementary Figure 1K-L), the false minimum is slightly more compact and has a more
9 ordered secondary structure. Two beta sheet regions in the false minimum are beta
10 strands/unordered in the native structure and two alpha helices in the false minimum are beta
11 strands in the native structure (Supplementary Figure 1F-G).

12 The case of 1SEN, which has false minima for both Rosetta and Amber, is similar to 1T2I in
13 that the false minima are more ordered than the native structure, although the native structure
14 forms more contacts between remote regions than do the false minima (Supplementary
15 Figure 1A-D). Residues 85-96 form a tight beta hairpin in the false minima, whereas the native
16 residues 85-96 has a longer loop between the beta strands, resulting in a shorter, less tight, beta
17 hairpin. Additionally, residues 94-109 in the native are entirely disordered, while that of the
18 false minimum begins as a beta strand and ends in an alpha helix (Supplementary Figure 1H).
19 Decoys that are predicted as false minima often have the same overall structure and contact maps
20 as native structures, yet secondary structure differences may result in large structural deviation.
21 In some cases, false minima contain more ordered secondary structures yet fewer contacts than
22 native conformations; the propensity away from disordered loops may result in lower energies
23 for these false minima.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

Per-residue Rosetta energy decomposition

To investigate whether certain residues, structural elements, or energy terms contribute more to false minima conformations, we analyzed the per-residue score decomposition for Rosetta scores for the three systems outlined above (2QY7, 1T2I, and 1SEN). We were unable to perform the same decomposition for Amber as the GB solvation term is not pairwise-decomposable³⁸. We calculated the Z-scores for each residue over the lowest-scoring native and false minimum conformations. We identified residues as possibly implicated in false minima if the false minimum residue Z-score score was lower than the native residue Z-score by at least one (i.e. the distance between the two was greater than one standard deviation). We have highlighted these residues (Figure 3A-C). False minima contributing residues were distributed over the conformations and did not cluster to any particular region. Moreover, false minima contributing residues were found in various types of secondary structure: alpha helices, beta strands, and loops. It is therefore currently not possible to attribute Rosetta false minima to any single per-residue propensity, but as expected, several small errors in energy estimation may lead to the observed incorrect scoring.



1
 2 **Figure 3.** Per-residue and per-score-term propensity of score-functions toward false minima. (A-
 3 C) Native (gray) and Rosetta-minimized (salmon) structures of 2QY7, 1T2I, and 1SEN
 4 respectively. Rosetta-minimized residues that are scored by Rosetta as greater than 1 standard
 5 deviation away from the corresponding native residue are highlighted in red. Heatmaps of per-

1 structure, score-term contribution to Rosetta-determined (D) and Amber-determined (E) false
2 minima and true maxima. The row marked Overall shows the percentage of structures that
3 indicate some degree of implication for that score-term.

4 **Per-scoreterm contributions of Amber and Rosetta**

5 We reasoned that insight about the performance and pathologies of each energy function could
6 be gained by identifying the energy terms that are responsible for correct and incorrect
7 evaluations within the same energy function. For example, we asked which terms in the Amber
8 energy function help it avoid mis-scoring a decoy (called Amber true maximum) that is
9 identified as a false minimum in the Rosetta landscape (called Rosetta false minimum), and vice
10 versa.

11 We identified terms that contribute to false minima and true maxima by calculating the Z-
12 scores per decoy set and native set for each protein. If the lowest native score-term Z-score is
13 greater than the false minimum score-term by at least one, that term is implicated in that false
14 minimum. The reverse (i.e. true maximum score-term Z-score is greater than the lowest native
15 score-term Z-score by at least one) is true for identifying true maximum contributing score-
16 terms. The heatmap in Figure 3D depicts the fraction of Rosetta false minima decoys (top) and
17 true maxima decoys (bottom) that show some degree of implication for each score-term. This is
18 calculated both on a per-protein basis and over the entire false minima/true maxima sets. Several
19 score-terms, including hbond_sr_bb, fa_dun, fa_rep and omega, are implicated in a majority of
20 false minima in the Rosetta talaris2014 energy function. A set of other score-terms contribute to
21 a majority of Rosetta true maxima (or Amber false minima), including rama, hbond_bb_sc,
22 hbond_sc, p_aa_pp, and fa_elec. These are score-terms that are not usually implicated in Rosetta
23 false minima, thus demonstrating that the score-terms that contribute to the two trends (towards

1 false minima and true maxima) are mutually exclusive. Except fa_elec, the other terms identified
2 as helping “rescue” Amber false minima are all PDB-statistics derived, and it is not surprising
3 that they are implicated in correcting the errors of the more physics-based Amber energy
4 function.

5 We next performed a similar analysis on Amber score-terms for both Amber false minima and
6 Amber true maxima (Figure 3E). We found that bond, angle, and gb are responsible for more
7 than 50% of Amber false minima and that dihedral and elec are implicated in rescuing Rosetta
8 false minima (Amber true maxima). We found that score-terms that are responsible for false
9 minima are not implicated in true maxima and vice versa. Similar to the identification of
10 statistically-derived terms in Rosetta as being responsible for correctly scoring Amber false
11 minima, we find that physics-based terms, i.e., elec (which is counterbalanced by gb) and
12 dihedral potentials, that are orthogonal to the talaris2014 Rosetta scorefunction, are implicated in
13 the rescue of Rosetta false minima by Amber.

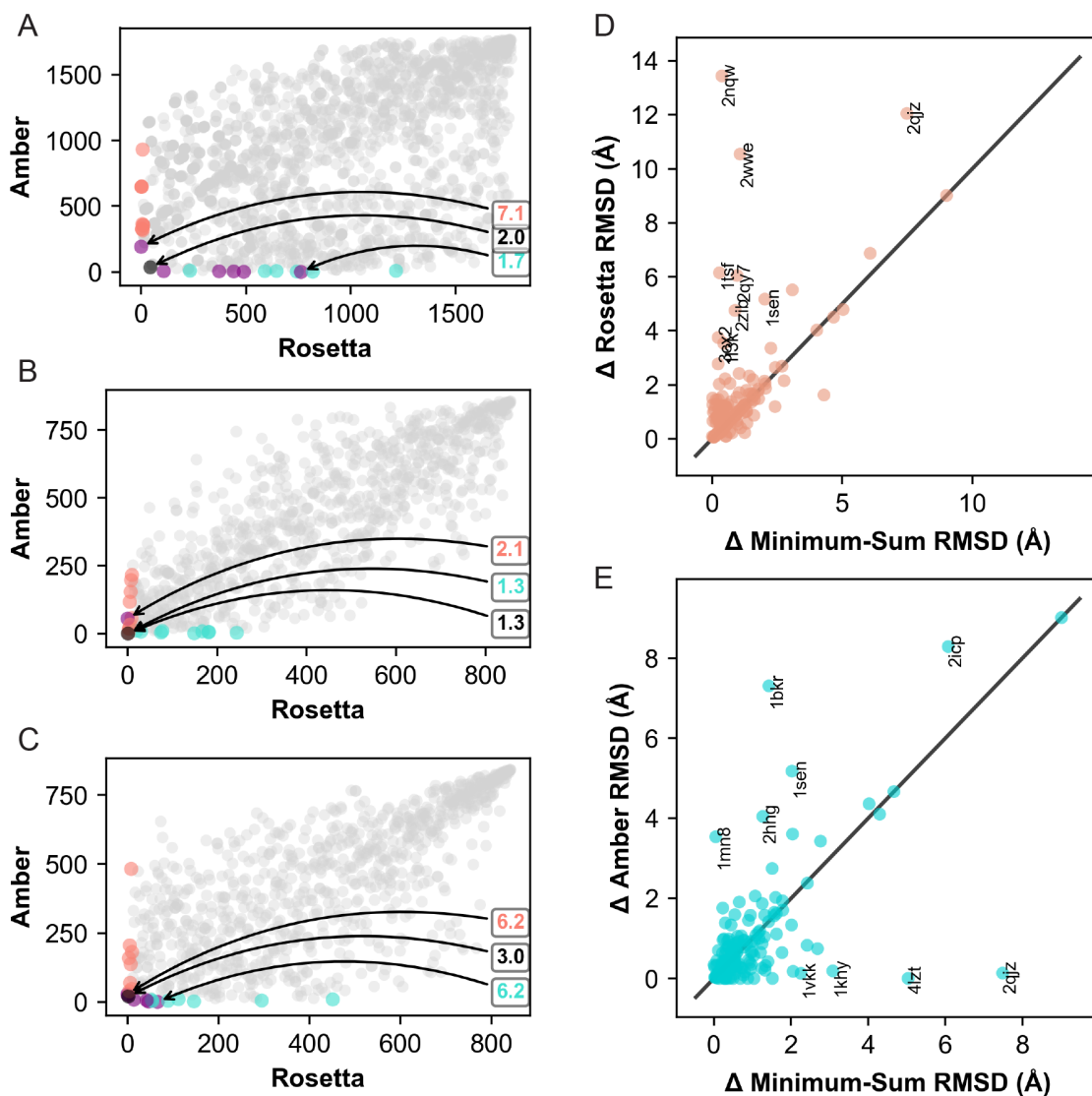
14

15 **Combining rankings to select decoys improves decoy selection**

16 Based on the results above indicating that the rescue of false minima in the landscape
17 generated by one energy function can be effected by the use of the other energy function due to
18 additional terms or different parameterization of terms, we sought to develop an approach to
19 productively combine the two landscapes for model selection. In model selection (for example in
20 protein structure prediction) the challenge is to select a near-native conformation from a set of
21 decoy conformations based on one or more energy values or other features. Typically, an energy
22 value obtained from a single energy function is used. In the current benchmark set, if model
23 selection is performed by the Rosetta and Amber energy functions individually, the Rosetta

1 lowest-scored decoy has an RMSD of $> 5.0 \text{ \AA}$ for thirteen out of 150 systems, while the lowest-
2 scored Amber decoy has an RMSD of $> 5.0 \text{ \AA}$ for seven systems (four of which overlap with the
3 aforementioned Rosetta systems). We designed a minimum-sum based algorithm (see Methods)
4 to select a decoy conformation based on both sets of ranks to improve the chances of selecting a
5 near-native decoy.

6 We found that our minimum-sum algorithm improved model selection for both Rosetta and
7 Amber rankings (Figure 4D-E), although it improved model selection for Rosetta to a greater
8 extent. The minimum-sum selected decoy had a lower RMSD than the lowest-scoring Rosetta
9 decoy by at least 1 \AA for ten out of the thirteen cases mentioned above and a lower RMSD than
10 the lowest-scoring Amber decoy by at least 1 \AA for four out of the seven cases mentioned above.
11 More generally, the minimum-sum selected decoy had a lower RMSD than the lowest-scoring
12 Rosetta decoy for 22 out of 150 cases and a lower RMSD than the lowest-scoring Amber decoy
13 for 11 out of 150 cases.



1
 2 **Figure 4.** Minimum-sum model selection. (A-C) Scatterplots of Rosetta-rank vs. Amber-rank for
 3 all decoys of 2QY7, 1T2I, and 1SEN respectively. Each point represents one decoy
 4 conformation. The set of Pareto solutions is purple, the top-10 ranked Amber decoys are
 5 turquoise, the top-10 ranked Rosetta decoys are salmon, and the minimum-sum solution is black.
 6 Annotations represent the RMSD in Å from native for the top-ranked Amber decoy (turquoise),
 7 top-ranked Rosetta decoy (salmon), and minimum-sum solution (black). Scatterplots that show
 8 the efficacy of the minimum-sum solution at minimizing the distance from native relative to the

1 top-ranked Rosetta decoy (D) and top-ranked Amber decoy (E). Each point represents one
2 system. The x-axis is the difference between the minimum-sum solution RMSD from native and
3 the RMSD of the minimum-RMSD decoy conformation, while the y-axis is the difference
4 between the Rosetta lowest-ranked conformation (D) or Amber lowest-ranked conformation (E)
5 RMSD from native and the RMSD of the minimum-RMSD decoy conformation. Points that fall
6 outside the 95% prediction interval are annotated.

7 We examined the false minima cases described above (2QY7, 1T2I, and 1SEN) and found that
8 the minimum-sum decoy generally had a lower RMSD than that of Rosetta- or Amber-selected
9 decoys (Figure 4A-C). However, for 2QY7, which contains a false minimum for Rosetta but not
10 for Amber, the Amber-selected decoy had a slightly lower RMSD than that of the minimum-sum
11 decoy (1.7 Å vs. 2.0 Å). Nevertheless, the minimum-sum selected decoy RMSD is significantly
12 lower than that of the Rosetta-selected decoy (2.0 Å vs. 7.1 Å). Thus, a minimum-sum
13 framework allows combining the two energy functions productively to select a near-native
14 model.

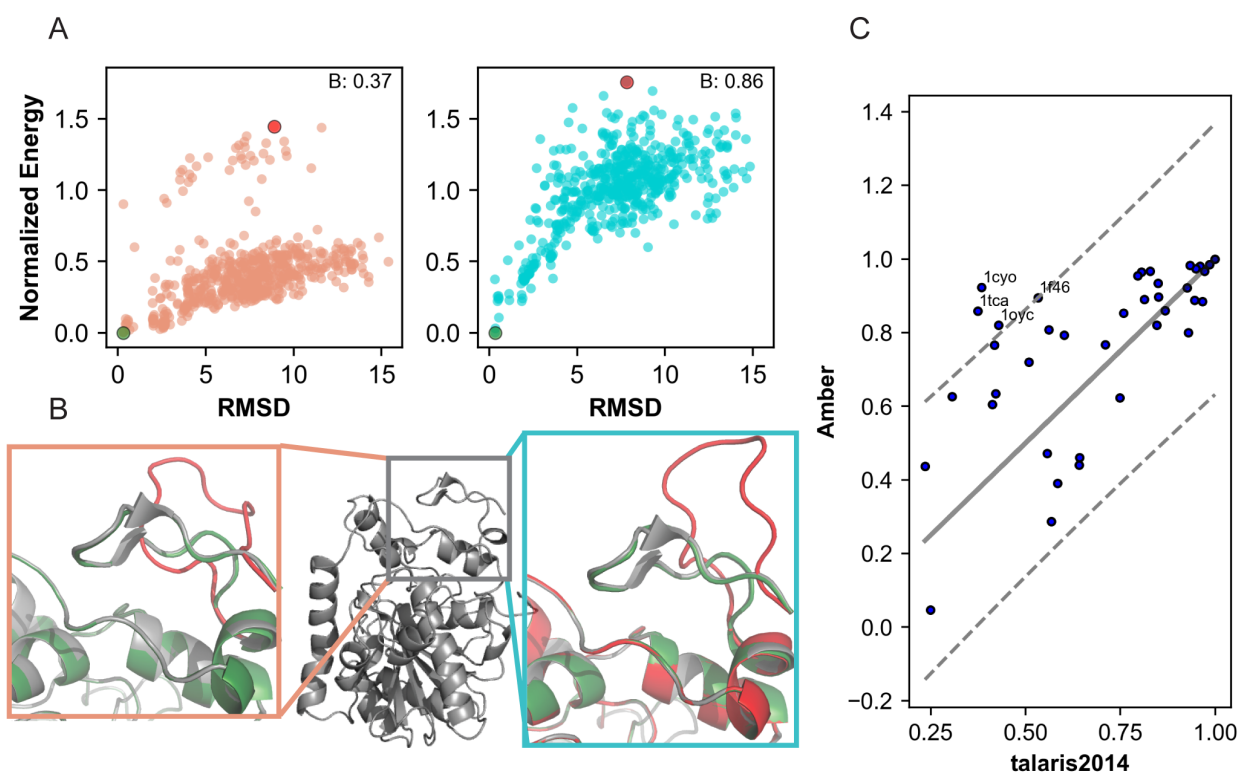
15

16 **Loop Modeling**

17 The conformational variability of loops plays a multi-functional role in protein structure and
18 function. They are implicated in stability and folding pathways³⁴, binding and active sites^{35,36}, and
19 binding other proteins^{37,38}. Efficient sampling algorithms have been developed³⁶⁻³⁹, but loop
20 structure prediction efforts can be limited by energy functions, as the energy gaps between loops
21 are smaller and minima are narrow⁴⁰. Therefore, we tested both Amber and Rosetta energy
22 functions on a loop modeling benchmark obtained from T. Kortemme and S. O'Connor. In this
23 benchmark, most of the structure remains the same over the set of decoys; the difference lies in a

1 small loop region, which can vary highly in RMSD. The energy gaps between structures are
2 therefore smaller; thus, loop modeling provides a more stringent test to distinguish between
3 energy functions.

4 We found that Amber ranked loops more accurately than did Rosetta (Figure 5C). Several
5 systems had significantly higher B values with Amber than with Rosetta. Figure 5A depicts the
6 energy landscapes for one of these structures (1TCA). The Amber funnel is steeper than that of
7 Rosetta, which is reflected in its higher B (0.86 vs. 0.37). The lowest-energy and highest-energy
8 loop conformations are shown for both Rosetta and Amber in Figure 5B. Both Rosetta and
9 Amber rank the lowest-energy and highest-energy conformations correctly.



10
11 **Figure 5.** Loop modeling benchmark. (A) Energy landscape for 1TCA. Each dot on the plot
12 represents one decoy conformation. The x-axis is RMSD from native and the y-axis is
13 normalized energy. The B metric, which represents the efficacy of the score-function at

1 differentiating between native and non-native decoys, is shown at the top right corner of each
2 plot. Rosetta plots are to the left, in salmon, and Amber plots are to the right, in turquoise. The
3 lowest-energy decoy conformation in each plot is shown in green and the highest-energy decoy
4 conformation is shown in red. (B) Native structure of 1TCA (gray) and close-ups of loop
5 conformations for lowest-energy decoys (green) and highest-energy decoys (red) for Rosetta
6 (salmon box) and Amber (turquoise box). (C) General performance of Rosetta talaris2014
7 scoring function vs. Amber scoring function over the entire loop modeling set. Each dot
8 represents the *B* metric for one system. The black line is $x=y$ and the dashed line represents the
9 95% prediction interval. Any points that lie outside the 95% prediction interval are annotated
10 with the PDB ID of that system.

11 **Discussion**

12 Systematic comparison of Amber ff14SBOonlySC (a physically-derived energy function) and
13 Rosetta talaris2014 (both physical and statistical based) reveals the strengths and weaknesses of
14 each energy function. Generally, Amber ff14SBOonlySC performs better than Rosetta talaris2014
15 at both decoy discrimination and loop modeling. However, comparison of Amber ff14SBOonlySC
16 to Rosetta REF2015 (the newer, default Rosetta energy function) reveals that REF2015, which
17 has more physically-derived terms than talaris2014, performs comparably well to Amber
18 ff14SBOonlySC. Examination of Rosetta talaris2014 score-terms that rescue Amber
19 ff14SBOonlySC false minima and Amber ff14SBOonlySC score-terms that correct Rosetta
20 talaris2014 false minima reveals two possible sources for the performance improvement of
21 REF2015. While two of the Rosetta score-terms and two of the Amber score-terms that
22 contribute to the correction of false minima are counterparts to each other (Amber dihedral and
23 Rosetta rama, and Amber elec and Rosetta fa_elec), subtle nuances in their derivation and

1 parameterization appear to influence the propensity of each energy function toward false
2 minima. Although rama and dihedral both score the propensity of the backbone dihedral angles,
3 rama does so in a statistically-derived manner while dihedral is based on fits to quantum
4 chemistry data. Both elec and fa_elec are derived from a Coulombic model, yet they are
5 differently parameterized; the Amber elec is parameterized *via* small-molecule properties,
6 whereas fa_elec is optimized on larger biomolecular structures. The improvement of Rosetta
7 REF2015 over Rosetta talaris2014 may be caused by its greater inclusion of physical-derived
8 terms (bond, angle, etc.) and/or its parameterization on both small-molecule properties and larger
9 biomolecular structures.

10 Model selection, or the ability to select a near-native decoy from a set of decoy conformations
11 is a general problem in protein structure prediction. If low-energy decoys exist in false minima
12 in the energy landscape, it is difficult to identify conformations that are near-native. Since
13 Amber and Rosetta provide different, semi-orthogonal information, a framework to combine the
14 two rankings enable the identification of near-native decoys. The minimum-sum based
15 algorithm that we have implemented improves model selection for 15% of structures over
16 Rosetta model selection and 7.3% of structures over Amber model selection. The model
17 selection algorithm is extensible to any two sets of energy functions or model ranks for one set of
18 models and can thus be used to combine any two sources of information to produce meaningful
19 improvements in near-native decoy selection.

20 The approach described here should enable comparative analysis and combination of future
21 versions of both Amber and Rosetta scoring functions, and enable a variety of biomolecular
22 modeling tasks.

23

1 ASSOCIATED CONTENT

2 **Supporting Information.**

3 The following files are available free of charge.

4 LoopDefs: definitions for the loops in the loop modeling benchmark (xlsx)

5 Supplementary_Software: descriptions of software and scripts used in Methods (docx)

6 Figures: S1, False minima contact maps and structures; S2, Energy landscapes for all decoy
7 discrimination systems; S3, plot of minimum, minimum-sum, Rosetta, and Amber RMSDs; S4,

8 Energy landscapes for all loop modeling systems. Tables: S1, B-metric values for Amber and

9 Rosetta decoy discrimination systems; S2, values of RMSD for minimum-sum-selected, Rosetta-

10 selected, and Amber-selected models as well as the actual lowest-RMSD value; S3, B-metric

11 values for Amber and Rosetta loop modeling systems (PDF)

12 AUTHOR INFORMATION

13 **Corresponding Author**

14 *Emails: sagar.khare@rutgers.edu, david.case@rutgers.edu

15 **Present Addresses**

16 †Schrodinger, Inc. 120 West 45th Street, 17th Floor, Tower 45, New York, NY 10036

17 **Author Contributions**

18 ABR, KMB, SDK and DAC wrote the manuscript. All authors have given approval to the final
19 version of the manuscript.

20 **Funding Sources**

1 This material is based upon work supported by RosettaCommons and the National Science
2 Foundation Graduate Research Fellowship under Grant No. DGE-1433187 (ABR).

3 ACKNOWLEDGMENT

4 We thank F. DiMaio, H. Park for the REF2015 dataset and S. O'Connor and T. Kortemme for
5 the loop modeling dataset.

6 ABBREVIATIONS

7 RMSD, root-mean-square-deviation; REF2015, Rosetta Energy Function 2015; LBFGS,
8 Limited-memory Broyden-Fletcher-Goldfarb-Shanno; lbfgs_armijo_nonmonotone, LBFGS
9 minimizer implementation with inexact line search conditions; GB-neck2, generalized Born
10 implicit solvent model.

11

12 REFERENCES

- 13 (1) Anfinsen, C. B. *Science* **1973**, *181* (4096), 223–230.
- 14 (2) Lu, H.; Skolnick, J. *Proteins Struct. Funct. Genet.* **2001**, *44* (3), 223–232.
- 15 (3) Brooks, B. R.; Brooks, C. L.; Mackerell, A. D.; Nilsson, L.; Petrella, R. J.; Roux, B.;
16 Won, Y.; Archontis, G.; Bartels, C.; Boresch, S.; Caflisch, A.; Caves, L.; Cui, Q.; Dinner,
17 A. R.; Feig, M.; Fischer, S.; Gao, J.; Hodoscek, M.; Im, W.; Kuczera, K.; Lazaridis, T.;
18 Ma, J.; Ovchinnikov, V.; Paci, E.; Pastor, R. W.; Post, C. B.; Pu, J. Z.; Schaefer, M.;
19 Tidor, B.; Venable, R. M.; Woodcock, H. L.; Wu, X.; Yang, W.; York, D. M.; Karplus,
20 M. *J. Comput. Chem.* **2009**, *30* (10), 1545–1614.
- 21 (4) Cornell, W. D.; Cieplak, P.; Bayly, C. I.; Gould, I. R.; Merz, K. M.; Ferguson, D. M.;

- 1 Spellmeyer, D. C.; Fox, T.; Caldwell, J. W.; Kollman, P. A. *J. Am. Chem. Soc.* **1995**, *117*
2 (19), 5179–5197.
- 3 (5) Jernigan, R. L.; Bahar, I. *Current Opinion in Structural Biology*. 1996, pp 195–209.
- 4 (6) Shen, M.-Y.; Sali, A. *Protein Sci.* **2006**, *15* (11), 2507–2524.
- 5 (7) Ponder, J. W.; Case, D. A. *Advances in Protein Chemistry*. 2003, pp 27–85.
- 6 (8) Simons, K. T.; Kooperberg, C.; Huang, E.; Baker, D. *J. Mol. Biol.* **1997**, *268* (1), 209–
7 225.
- 8 (9) Jorgensen, W. L.; Maxwell, D. S.; Tirado-Rives, J. *J. Am. Chem. Soc.* **1996**, *118* (45),
9 11225–11236.
- 10 (10) Maier, J. A.; Martinez, C.; Kasavajhala, K.; Wickstrom, L.; Hauser, K. E.; Simmerling, C.
11 *J. Chem. Theory Comput.* **2015**, *11* (8), 3696–3713.
- 12 (11) MacKerell, A. D.; Bashford, D.; Bellott, M.; Dunbrack, R. L.; Evanseck, J. D.; Field, M.
13 J.; Fischer, S.; Gao, J.; Guo, H.; Ha, S.; Joseph-McCarthy, D.; Kuchnir, L.; Kuczera, K.;
14 Lau, F. T. K.; Mattos, C.; Michnick, S.; Ngo, T.; Nguyen, D. T.; Prodhom, B.; Reiher, W.
15 E.; Roux, B.; Schlenkrich, M.; Smith, J. C.; Stote, R.; Straub, J.; Watanabe, M.;
16 Wiórkiewicz-Kuczera, J.; Yin, D.; Karplus, M. *J. Phys. Chem. B* **1998**, *102* (18), 3586–
17 3616.
- 18 (12) Xu, D.; Zhang, Y. *Proteins Struct. Funct. Bioinforma.* **2012**, *80* (7), 1715–1735.
- 19 (13) O’Meara, M. J.; Leaver-Fay, A.; Tyka, M. D.; Stein, A.; Houlihan, K.; Dimaio, F.;
20 Bradley, P.; Kortemme, T.; Baker, D.; Snoeyink, J.; Kuhlman, B. *J. Chem. Theory*

- 1 *Comput.* **2015**, *11* (2), 609–622.
- 2 (14) Shapovalov, M. V.; Dunbrack, R. L. *Structure* **2011**, *19* (6), 844–858.
- 3 (15) Best, R. B.; Zhu, X.; Shim, J.; Lopes, P. E. M.; Mittal, J.; Feig, M.; MacKerell, A. D. *J.*
4 *Chem. Theory Comput.* **2012**, *8* (9), 3257–3273.
- 5 (16) Raval, A.; Piana, S.; Eastwood, M. P.; Dror, R. O.; Shaw, D. E. *Proteins Struct. Funct.*
6 *Bioinforma.* **2012**, *80* (8), 2071–2079.
- 7 (17) Lindorff-Larsen, K.; Maragakis, P.; Piana, S.; Eastwood, M. P.; Dror, R. O.; Shaw, D. E.
8 *PLoS One* **2012**, *7* (2).
- 9 (18) Park, H.; Bradley, P.; Greisen, P.; Liu, Y.; Mulligan, V. K.; Kim, D. E.; Baker, D.;
10 Dimaio, F. J. *J. Chem. Theory Comput.* **2016**, *12* (12), 6201–6212.
- 11 (19) Alford, R. F.; Leaver-Fay, A.; Jeliazkov, J. R.; O’Meara, M. J.; DiMaio, F. P.; Park, H.;
12 Shapovalov, M. V.; Renfrew, P. D.; Mulligan, V. K.; Kappel, K.; Labonte, J. W.; Pacella,
13 M. S.; Bonneau, R.; Bradley, P.; Dunbrack, R. L.; Das, R.; Baker, D.; Kuhlman, B.;
14 Kortemme, T.; Gray, J. J. *J. Chem. Theory Comput.* **2017**, *13* (6), 3031–3048.
- 15 (20) Beauchamp, K. A.; Lin, Y. S.; Das, R.; Pande, V. S. *J. Chem. Theory Comput.* **2012**, *8* (4),
16 1409–1414.
- 17 (21) Cino, E. A.; Choy, W. Y.; Karttunen, M. *J. Chem. Theory Comput.* **2012**, *8* (8), 2725–
18 2740.
- 19 (22) Nguyen, H.; Roe, D. R.; Simmerling, C. *J. Chem. Theory Comput.* **2013**, *9* (4), 2020–
20 2034.

- 1 (23) Nguyen, H.; Maier, J.; Huang, H.; Perrone, V.; Simmerling, C. *J. Am. Chem. Soc.* **2014**,
2 *136* (40), 13959–13962.
- 3 (24) Conway, P.; Tyka, M. D.; DiMaio, F.; Konerding, D. E.; Baker, D. *Protein Sci.* **2014**, *23*
4 (1), 47–55.
- 5 (25) Berman, H. M.; Westbrook, J.; Feng, Z.; Gilliland, G.; Bhat, T. N.; Weissig, H.;
6 Shindyalov, I. N.; Bourne, P. E. *Nucleic Acids Res.* **2000**, *28*, 235–242.
- 7 (26) Tyka, M. D.; Keedy, D. A.; Andre, I.; Dimaio, F.; Song, Y.; Richardson, D. C.;
8 Richardson, J. S.; Baker, D. *J. Mol. Biol.* **2011**, *405* (2), 607–618.
- 9 (27) Tyka, M. D.; Jung, K.; Baker, D. *J. Comput. Chem.* **2012**, *33* (31), 2483–2491.
- 10 (28) Case, D. A.; Betz, R. M.; Cerutti, D. S.; Cheatham III, T. E.; Darden, T. A.; Duke, R. E.;
11 Giese, T. J.; Gohlke, H.; Goetz, A. W.; Homeyer, N.; Izadi, S.; Janowski, P.; Kaus, J.;
12 Kovalenko, A.; Lee, T. S.; LeGrand, S.; Li, P.; Lin, C.; Luchko, T.; Luo, R.; Madej, B.;
13 Mermelstein, D.; Merz, K. M.; Monard, G.; Nguyen, H.; Nguyen, H. T.; Omelyan, I.;
14 Onufriev, A.; Roe, D. R.; Roitberg, A.; Sagui, C.; Simmerling, C. L.; Botello-Smith, W.
15 M.; Swails, J.; Walker, R. C.; Wang, J.; Wolf, R. M.; Wu, X.; Xiao, L.; Kollman, P. A.
16 *AMBER 2016*; University of California: San Francisco, 2016.
- 17 (29) Liu, D. C.; Nocedal, J. *Math. Program.* **1989**, *45* (1–3), 503–528.
- 18 (30) Nguyen, H.; Roe, D. R.; Swails, J. M.; Case, D. A. *Manuscr. Prep.* **2017**.
- 19 (31) Dill, K. a; MacCallum, J. L. *Science* **2012**, *338* (6110), 1042–1046.
- 20 (32) Wolynes, P. G.; Onuchic, J. N.; Thirumalai, D. *Science* **1995**, *267* (5204), 1619–1620.

- 1 (33) Shakhnovich, E. *Chemical Reviews*. 2006, pp 1559–1588.
- 2 (34) Xiang, Z.; Soto, C. S.; Honig, B. *Proc. Natl. Acad. Sci. U. S. A.* **2002**, *99*, 7432–7437.
- 3 (35) Fiser, A.; Kinh Gian Do, R.; Sali. *Protein Sci.* **2000**, *9*, 1753–1773.
- 4 (36) Murphy, P. M.; Bolduc, J. M.; Gallaher, J. L.; Stoddard, B. L.; Baker, D. *Proc. Natl.*
5 *Acad. Sci. U. S. A.* **2009**, *106* (23), 9215–9220.
- 6 (37) Wang, C.; Bradley, P.; Baker, D. *J. Mol. Biol.* **2007**, *373* (2), 503–519.
- 7 (38) Mandell, D. J.; Coutsiar, E. a; Kortemme, T. *Nat. Methods* **2009**, *6* (8), 551–552.
- 8 (39) Ollikainen, N.; Smith, C. A.; Fraser, J. S.; Kortemme, T. *Methods Enzymol.* **2013**, *18* (9),
9 1199–1216.
- 10 (40) Stein, A.; Kortemme, T. *PLoS One* **2013**, *8* (5).

11

12

13

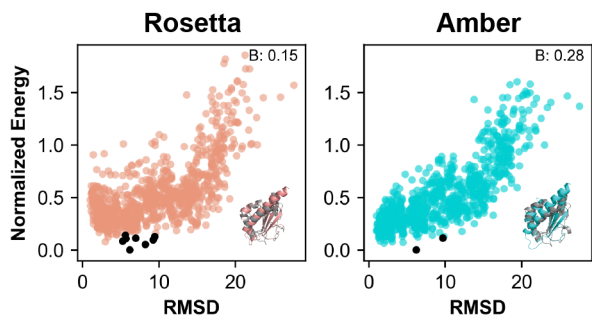
14

15

16

17

18 For Table of Contents Only



1