

EM-Fold: De Novo Atomic-Detail Protein Structure Determination from Medium-Resolution Density Maps

Steffen Lindert,¹ Nathan Alexander,¹ Nils Wötzel,¹ Mert Karakaş,¹ Phoebe L. Stewart,² and Jens Meiler^{1,*}

¹Department of Chemistry and Center for Structural Biology, Vanderbilt University, Nashville, TN 37212, USA

²Department of Pharmacology and Cleveland Center for Membrane and Structural Biology, Case Western Reserve University, Cleveland, OH 44106, USA

*Correspondence: jens.meiler@vanderbilt.edu

DOI 10.1016/j.str.2012.01.023

SUMMARY

Electron density maps of membrane proteins or large macromolecular complexes are frequently only determined at medium resolution between 4 Å and 10 Å, either by cryo-electron microscopy or X-ray crystallography. In these density maps, the general arrangement of secondary structure elements (SSEs) is revealed, whereas their directionality and connectivity remain elusive. We demonstrate that the topology of proteins with up to 250 amino acids can be determined from such density maps when combined with a computational protein folding protocol. Furthermore, we accurately reconstruct atomic detail in loop regions and amino acid side chains not visible in the experimental data. The EM-Fold algorithm assembles the SSEs de novo before atomic detail is added using Rosetta. In a benchmark of 27 proteins, the protocol consistently and reproducibly achieves models with root mean square deviation values <3 Å.

INTRODUCTION

In the field of protein structure determination cryo-electron microscopy (cryoEM) has been established as a viable approach for studying the structure and dynamics of macromolecular structure of large protein complexes at near native conditions. CryoEM is invaluable in cases where alternative approaches such as X-ray crystallography and nuclear magnetic resonance (NMR) fail. In recent years, cryoEM density maps have reached high enough resolutions to provide sufficient detail to trace the protein backbone (Liu et al., 2010a; Ludtke et al., 2008; Zhang et al., 2011; Zhou, 2008). More routinely, resolutions <10 Å are reached that reveal the location of α helices (Cong et al., 2010; Liu et al., 2010b; Ludtke et al., 2008, 2004; Min et al., 2006; Saban et al., 2006; Serysheva et al., 2008; Villa et al., 2009). Additionally, β strands become visible at resolutions around 6 Å. However, connectivity and directionality of these secondary structure elements (SSEs) and their alignment with the primary protein sequence remains ambiguous in these medium-resolution density maps (5–10 Å resolution)—that is, it remains unknown which part of the protein's primary sequence forms which α helix or which β strand and where the N and C termini of the helices and strands reside. Computational methods are

needed to help resolve this ambiguity. Several algorithms that help identify SSEs in a density map have been published. α Helices and β strand regions can be identified automatically by methods using segmentation and feature extraction (Baker et al., 2007; Dal Palù et al., 2006; Jiang et al., 2001; Kong et al., 2004). Furthermore, even high-quality medium-resolution cryoEM maps typically lack information at atomic detail, such as the conformation of loops and side chains. We explore the potential of computational methods to aid in the interpretation of maps by reconstructing structural information that is not readily visible at the respective resolution.

In the past, numerous experimental techniques have been successfully combined with computational methods. A combination of computational algorithms with sparse structural information from NMR spectroscopy (Bowers et al., 2000; Meiler and Baker, 2003b, 2005; Rohl and Baker, 2002) and electron paramagnetic resonance spectroscopy (Alexander et al., 2008; Hanson et al., 2008; Hirst et al., 2011) experiments has led to the construction of protein models that are accurate at atomic detail. Final models include atomic detail that is beyond the resolution of the experiment because of judicious use of complementary computational algorithms. A prerequisite for success in this regard is that the experimental data restrain the conformational space sufficiently to allow sampling of protein backbone conformations at a distance of about 1–2 Å around the global energy minimum. As a result, some protein models will have root mean square deviations (RMSDs) from the correct structure of <3.0 Å when normalized to a 100 residue protein (RMSD100) (Carugo and Pongor, 2001). This level of accuracy is sufficient to construct side chain coordinates in the protein core and allows discrimination of incorrect protein models on the basis of inferior energy values (Bradley et al., 2005). Several methods, such as Rosetta, Modeler, and EM-IMO, can be applied to the refinement of comparative or hand-traced models guided by cryoEM density maps (DiMaio et al., 2009; Topf et al., 2006; Zhu et al., 2010).

Here we demonstrate de novo protein structure determination to a level with accurate atomic detail using medium-resolution density maps to restrain the simulation. Our protocol consists of two steps: (1) determination of protein topology with an improved version of EM-Fold (Lindert et al., 2009) and (2) refinement to atomic detail accuracy using Rosetta (DiMaio et al., 2009). These two methods are highly complementary. EM-Fold differs from Rosetta in that it is tailored toward efficient sampling of the conformational search space at the cost of somewhat lower precision in its scoring functions. Its ability to move entire

SSEs as one element makes it well suited for folding into medium-resolution density maps. EM-Fold also differs from other cryoEM map-based model building algorithms such as SSEhunter (Baker et al., 2007) in that it is essentially a de novo protein folding tool that uses the cryoEM map as a restraint. EM-Fold builds topological models for a protein of interest that agree with the density map and fulfill basic requirements of protein structure. These models contain only SSEs, no coordinates for loop regions, and no side chains. It was originally developed to build models of α -helical proteins into medium-resolution density maps. Here, we present an updated version of the program that can place both types of SSEs (α helices and β strands) into the density map. Additional improvements to the algorithm include better handling of incorrect secondary structure prediction as well as a more advanced refinement protocol. EM-Fold models provide a good starting point for the Rosetta electron density refinement that also constructs loops and side chains guided by the cryoEM density map. The performance of the new folding and refinement protocol was tested on a benchmark set of 20 α -helical and seven β sheet proteins, 13 of which could be refined to atomic resolution detail. If SSEs are visible in the maps, this protocol could also be applied to low-resolution X-ray crystallography maps such as those obtained recently for several important membrane proteins (Ward et al., 2007) and macromolecular assemblies (Sibanda et al., 2010).

RESULTS AND DISCUSSION

EM-Fold determines the topology of a protein through placement of predicted SSEs into a cryoEM density map (Lindert et al., 2009). It uses a Monte Carlo Metropolis algorithm with a knowledge-based energy function that builds and refines physically realistic models that agree with the density map. The models are constructed from a pool of predicted SSEs. Model changes applied during the folding simulation include addition and deletion of SSEs together with swaps and rotations of these elements. To achieve higher accuracy in the initial models and aid atomic detail refinement, EM-Fold was extended to allow bending, translation, and dynamic length resizing of SSEs. For the models to accurately reflect such detail, the scoring function was adapted for direct comparison of the models with the density map. Furthermore, the present protocol employed a recently added feature in Rosetta that allows construction of loops and side chains guided by a density map (DiMaio et al., 2009). This is critical as we have found that accurate construction of the protein backbone in loop regions that connect SSEs is crucial for successful atomic-level refinement. Rosetta systematically rebuilds regions of the protein backbone that agree the least with the density map.

Benchmark Database of 20 α -Helical and Seven β Sheet Proteins with 150–250 Residues

A benchmark set of 20 α -helical and seven β sheet proteins with 150–250 amino acids was chosen to test the algorithm. The benchmark was limited to proteins up to 250 residues, as this provided the desired range of successes/failures demonstrating the capabilities and limitations of the method. As algorithms advance and computers become faster, the benchmark should

be expanded to contain larger proteins. The benchmark set is primarily composed of α -helical proteins, since these represent the majority of application cases, as α helices are observed more readily than β strands at medium resolution. However, the performance was also tested on seven proteins with β sheets to demonstrate general applicability. Density maps at 7 Å resolution were simulated for the 20 α -helical proteins. Density maps at 5 Å resolution were simulated for the seven β sheet containing proteins. At these resolutions, α helices and β strands can be unambiguously identified through visual inspection. Since the maps were simulated at 5 Å (β sheet containing proteins) and 7 Å (α -helical proteins) resolution, respectively, these are considered the resolution limits of the EM-Fold method. Maps at higher resolution will likely perform at least as well with EM-Fold, as more features tend to be present in these density maps. For higher-resolution maps, however, it might be advantageous to use methods designed to trace the protein backbone if the resolution of the density map allows (Baker et al., 2011).

Results of EM-Fold Assembly Step with Perfect and Realistic Secondary Structure Prediction

A schematic representation of the stages of the benchmark is shown in Figure 1. The algorithm adds, deletes, swaps, and flips as well as dynamically grows and shrinks SSEs to account for inaccurate secondary structure prediction. To avoid formation of unlikely SSEs, this new move is accompanied by scoring the agreement of the model's secondary structure with predicted secondary structure. To assess the impact of inaccurate secondary structure prediction on the algorithm, the benchmark was performed in two stages: using perfect and realistic secondary structure prediction, respectively. Realistic secondary structure prediction was generated as a consensus prediction of the methods used without any manual adjustment. Assuming perfect secondary structure prediction, the true topology was found among the top 20 scoring topologies for all but one of the 27 proteins. The exception was 1WBA, for which the correct topology rank was 287 (see Table S1 available online). Correct topology was defined as having placed all SSEs into the density rods that correspond to the location and orientation of that SSE in the experimental structure. Table 1 displays the results of the assembly runs for all 27 proteins with realistic secondary structure prediction, where RMSD100 values were calculated over all backbone atoms. In the case of realistic secondary prediction, the true topology was constructed in 23 out of the 27 proteins in the benchmark. For 15 of the 20 α -helical proteins and four of the seven β sheet proteins, the correct topology also ranked among the top 150 scoring topologies. For four additional proteins (1Z3Y, 2FQ4, 2IU1, 2NR7), the correct topology was constructed, but it was not identified among the best 150 topologies. These results indicate that it is more difficult to predict the correct topology of β -sheet-containing proteins. This may be because of a higher curvature of β strands in β sheets compared to α helices or the generally higher contact order of β proteins. A total of 19 out of 27 protein topologies (70%) were identified within the best 150 scoring topologies after the assembly step. While a higher success rate was definitely desirable, a 70% success rate is very competitive compared to other de novo folding benchmarks.

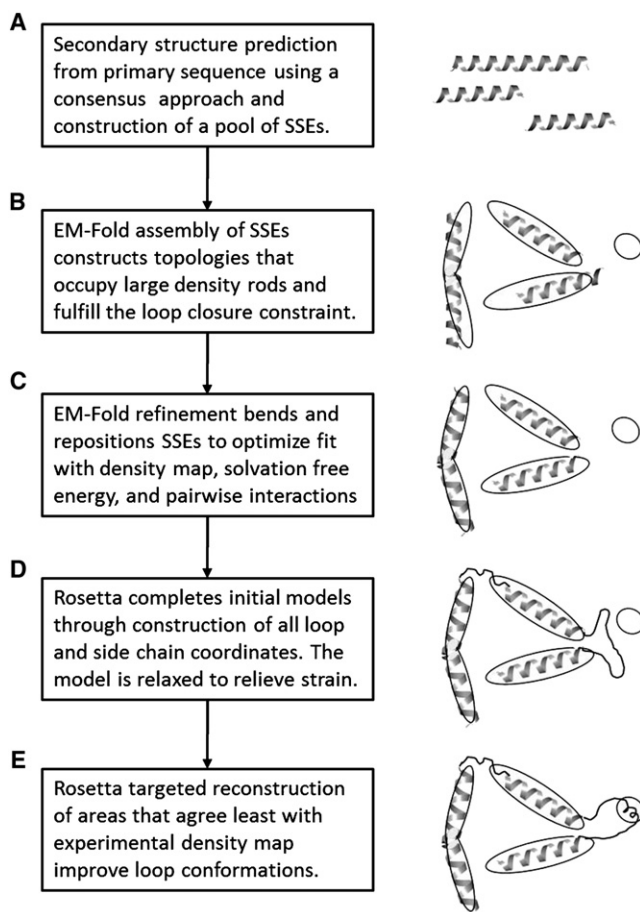


Figure 1. Schematic Representation of the Folding Protocol Used in the Benchmark

The scheme represents a three-density rod density map.

(A) Using a consensus of three secondary structure prediction methods, likely positions of long stretches of secondary structure in the primary sequence of the protein are identified. These positions are collected in a pool of idealized SSEs.

(B) The EM-Fold assembly step builds 50,000 models by assembling predicted, idealized SSEs into the identified density rods. The models contain no residues in the loop regions and no side chains. The top-scoring 150 topologies are carried over to the next step.

(C) The EM-Fold refinement step builds 500 refined models for each of the 150 topologies. The models generated in the assembly step are refined to better fit the density map. In particular, bending of idealized SSEs is performed. The top-scoring 50 topologies are carried over to the next step.

(D) Rosetta (round 1) builds loop models for each of the 50 topologies. Loops are built and the overall structure is relaxed. The top-scoring 15 topologies are carried over to the next step.

(E) Rosetta (round 2 and 3) identifies regions in the proteins that agree least with the density map and selectively rebuilds these identified regions and relaxes the entire structure. The top-scoring five topologies after round 2 are carried over into round 3.

EM-Fold Refinement Protocol Improves RMSD100s of Models

The top 150 scoring models selected after the assembly step were refined in EM-Fold with perturbations, including bending in addition to rotations and translations of the individual SSEs as described in Lindert et al. (2009). The agreement of the model

with the density map was scored using a cross-correlation coefficient. Table 1 shows the improvement in RMSD100 from the assembly step to the refinement step. RMSD100s were calculated over all backbone atoms. For all but one of the 19 successful proteins, the refinement step generated models that were lower in RMSD100 than the assembly step model. The maximal improvement was 2.4 Å for 1DVO (a model with RMSD100 of 1.32 Å was built but did not score best), and the average improvement was 1.0 Å. When only considering the best-scoring true topology models after refinement, the average improvement in RMSD100 was 0.2 Å, while the best improvement was 2.0 Å (for 1X91). For all but one protein, the correct topologies were among the top 50 scoring topologies after the refinement step. The exception was 1CHD, where the best-scoring topology rank was 87 after the refinement step. In particular, proteins where the true topology ranked worse than 25 after the assembly step were considerably improved in ranking by refinement. The top-scoring model for each of the 50 top-scoring topologies after the refinement step was used in the first round of the Rosetta refinement protocol.

Rosetta Refinement Improves Models and Reaches Atomic Detail Accuracy for Favorable Cases

An iterative refinement protocol was applied using Rosetta. The first round built loops and side chains for the 50 top scoring topologies from the EM-Fold refinement step. The resulting models underwent relaxation in the Rosetta force field. Regions that agreed least with the density map in the best-scoring 15 topologies were identified using Rosetta's `loops_from_density.linuxgccrelease` executable. These regions were rebuilt in a second round of the Rosetta refinement, followed by another relaxation of the models. Finally, the regions with the largest discrepancies to the density map in the top five scoring topologies after round 2 were rebuilt in round 3. Table 1 summarizes the results after each of the three rounds of Rosetta refinement. Fourteen of the 19 final best-scoring models corresponded to the correct topology. In the remaining cases, the true topology was ranked second in three cases and fourth in the two worst cases (Table S2 lists the RMSD100 values of the top-scoring models for completeness). Rosetta was thus able to identify the correct topology by score whenever a model with an RMSD100 <2.8 Å was built. This was the case for 14 of 19 proteins. The RMSD100s of the correct models after completion of the iterative refinement protocol ranged from 1.3 to 6.9 Å over the full length of the proteins and from 0.8 to 3.8 Å over the SSEs. The average RMSD100 was 3.0 Å over the full length of the protein and 2.2 Å over the SSEs. Thirteen of the proteins had backbone atom RMSD100s of <3.0 Å over all residues, indicating correct atomic detail accuracy. Figure 2 shows models for 1X91, 1OZ9, and 1DVO superimposed with the native structure, where RMSD100s of 1.1 Å, 1.4 Å, and 1.8 Å, respectively, over all residues were achieved. Side chain conformations in the protein core are shown for both the model and the native structure. The RMSD100 versus score plots for all three proteins are displayed next to the models. Figure S1 depicts the model evolution over all the rounds of the protocol for two of the proteins (1X91 and 1OZ9). It can be seen how the quality of the models improves from the EM-Fold assembly step to the third round of Rosetta refinement. Figures 3, 4, and

Table 1. Results of Benchmark on Set of 27 α , α/β , and β Proteins

Protein (PDB ID)	Size ^a	Rank/RMSD100, Å ^b		Rank/RMSD100 (RMSD100 SSEs), Å ^c			Rotamer Recovery ^d
		EM-Fold Assembly	EM-Fold Refinement	Rosetta Round 1	Rosetta Round 2	Rosetta Round 3	
α-Proteins							
1DVO	152, 4, 0	3/3.73	19/2.25	1/2.49 (1.30)	1/2.23 (1.33)	1/2.07 (1.36)	0.64
1GS9	165, 4, 0	31/3.83	15/4.01	3/3.76 (3.59)	5/3.87 (3.69)	4/3.96 (3.70)	0.43
1IAP	211, 7, 0	2/2.57	8/2.03	1/2.51 (1.31)	1/2.12 (1.27)	1/2.43 (1.28)	0.49
1ILK	151, 5, 0	73/3.08	23/3.17	1/2.78 (2.60)	1/2.75 (2.60)	1/2.64 (2.53)	1.00
1NIG	152, 4, 0	66/5.97	22/6.42	4/6.04 (4.44)	3/5.98 (4.47)	2/5.92 (4.37)	0.50
1OXJ	173, 4, 0	71/3.20	39/2.62	1/5.89 (1.50)	1/4.34 (1.60)	1/4.14 (1.68)	0.54
1X91	153, 5, 0	1/3.30	1/1.33	1/1.37 (0.87)	1/1.32 (0.92)	1/1.29 (0.78)	0.86
1Z3Y	238, 7, 0	621/-	-/-	-/-	-/-	-/-	-
2A6B	234, 6, 0	131/2.83	24/2.55	1/3.90 (1.76)	1/3.12 (1.74)	1/2.22 (1.80)	0.60
2FD5	180, 6, 0	1/2.88	37/2.09	1/1.76 (1.17)	1/1.68 (1.11)	1/2.19 (1.64)	0.54
2FM9	215, 9, 0	126/3.09	1/2.38	1/2.63 (2.29)	1/2.29 (2.11)	1/2.29 (2.02)	0.61
2FQ4	192, 7, 0	260/-	-/-	-/-	-/-	-/-	-
2G7S	194, 6, 0	1/2.47	29/2.52	1/2.04 (1.67)	1/2.00 (1.70)	1/2.01 (1.74)	0.58
2GEN	197, 7, 0	2/2.70	18/2.76	4/2.67 (2.25)	1/2.27 (2.05)	1/2.35 (2.08)	0.56
2IGC	164, 4, 0	23/4.51	41/6.23	1/7.10 (4.19)	5/6.91 (3.81)	2/6.93 (3.78)	0.53
2IOS	150, 6, 0	60/4.13	14/2.68	1/3.87 (3.03)	1/3.48 (3.18)	1/3.31 (3.04)	0.49
2IU1	208, 5, 0	849/-	-/-	-/-	-/-	-/-	-
2NR7	195, 5, 0	386/-	-/-	-/-	-/-	-/-	-
2O8P	227, 9, 0	15/2.82	18/2.77	2/2.35 (2.25)	2/2.05 (2.01)	1/2.18 (2.13)	0.48
2QK1	249, 9, 0	-/-	-/-	-/-	-/-	-/-	-
α/β Proteins							
1BJ7	156, 1, 8	-/-	-/-	-/-	-/-	-/-	-
1CHD	203, 1, 8	24/1.68	87/1.65	2/15.76 (1.5)	4/15.44 (1.49)	4/15.42 (1.5)	0.53
1ICX	155, 1, 7	131/2.40	47/3.84	1/2.51 (2.08)	1/2.35 (1.89)	1/2.17 (1.76)	0.57
1JL1	155, 3, 5	32/3.10	13/3.56	1/3.38 (2.85)	1/2.84 (2.03)	2/2.91 (2.30)	0.71
1OZ9	150, 5, 4	1/2.27	9/1.67	1/1.88 (1.21)	1/1.89 (1.41)	1/2.19 (1.85)	0.55
β Proteins							
1WBA	175, 0, 10	-/-	-/-	-/-	-/-	-/-	-
2QVK	192, 0, 7	-/-	-/-	-/-	-/-	-/-	-
Average RMSD100s		3.19	2.98	3.27 (2.20)	2.97 (2.13)	2.96 (2.18)	

All RMSD100 values are determined over the backbone atoms N, C _{α} , C and O. The proteins from the benchmark set that are considered a success after the EM-Fold assembly and refinement steps as well as after the third round of Rosetta refinement are shown in bold. The criteria for the individual success assignments were: correct topology within the top 150 scoring models after the EM-Fold assembly step, correct topology within the top 50 scoring models after the EM-Fold refinement step, and correct topology being the top-scoring models with an RMSD100 of <3 Å after the third round of Rosetta refinement.

See also [Tables S1, S2, and S3](#).

^aNumber of amino acids, number of α helices with at least 12 residues, and number of β strands with at least five residues. Realistic secondary structure prediction has been used for the benchmark.

^bThe rank of the correct topology model within all scored models as well as the RMSD100 of the correct topology model are given.

^cEach column lists the rank of the correct topology model within all scored models, the RMSD100 of the correct topology model, and the RMSD100 of the correct topology model over residues in SSEs (numbers in parentheses).

^dRecovery in protein core is defined as the model having the same rotamers for all side chain angles. Core of the protein is defined as at least 22 neighbors.

S2 display score versus RMSD100 plots and the best models for all benchmark proteins. A clear funnel shape is visible for 14 out of 19 benchmark cases, with models having low RMSD100 values scoring better than models with high RMSD100 values. Occasionally, models with higher RMSD100 had scores that approached the score of the best-scoring

models with the correct topology (Figure 2C). An overlay of such a structure with the native model is shown in Figure 2D. All α helices were placed in the correct density rods, with one being placed in the wrong orientation. Most of the native helical interfaces were still present in this model, which explains its superior energy.

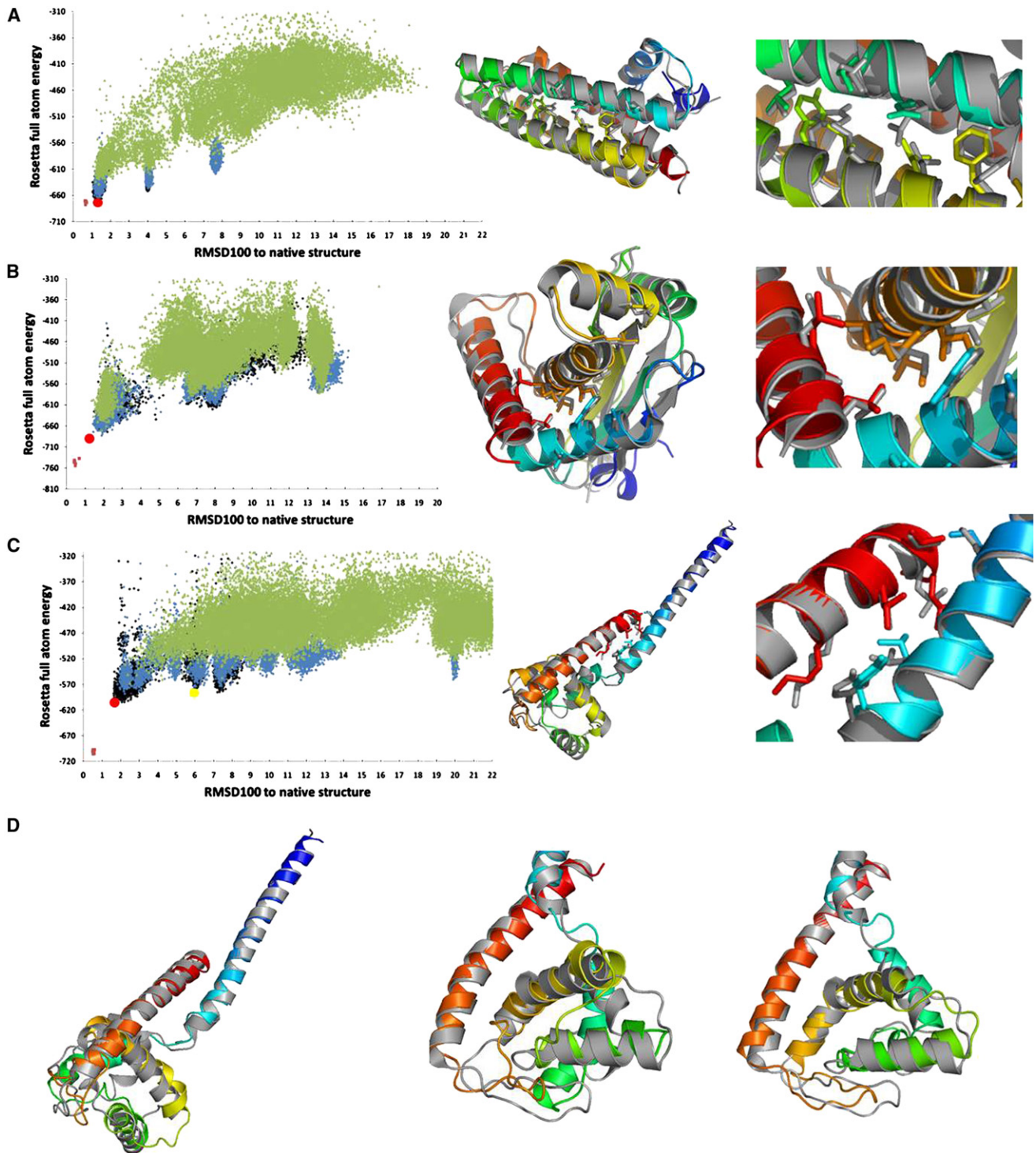


Figure 2. Rosetta Refinement RMSD100 versus Rosetta Energy Plots and Superimposition of Final Models after Rosetta Refinement with Medium-Sized Native Structures

(A–C) Energy plots for 1X91 (A), 1OZ9 (B), and 1DVO (C) are shown. Models from round 1 (green), round 2 (blue), and round 3 (black) of the Rosetta refinement are plotted. The native structure relaxed in Rosetta's force field is shown in violet for comparison. RMSD100s are calculated over all backbone atoms. For all three proteins, a model funnel is visible in the plot and the models corresponding to the correct topology score best allowing identification of the correct fold by score. Superimposition of the final models (rainbow-colored) of 1X91 (A), 1OZ9 (B), and 1DVO (C) with the original PDB structures (gray) are shown next to the energy funnels. The superimposed models are marked by a red dot in the energy funnels. A close-up view of side chain conformations in interfaces between SSEs is shown. (A) 1X91 has 153 residues. The model shown has an RMSD100 of 1.07 Å over the full length of the protein and 0.63 Å over the helical residues. (B) 1OZ9

The positive predictive value (PPV or precision) of the method to predict models with an RMSD100 below 3.0 Å was calculated after each of the three rounds of refinement. The PPV was 0.34 in round 1, 0.51 in round 2, and 0.50 in round 3. This indicates that there was a significant improvement in model quality and our ability to select good models by score when going from round 1 to 2. It also shows that the refinement process converged after round 2, with no further improvement when moving to round 3.

Quality Measure Can Distinguish between Successful and Unsuccessful Cases

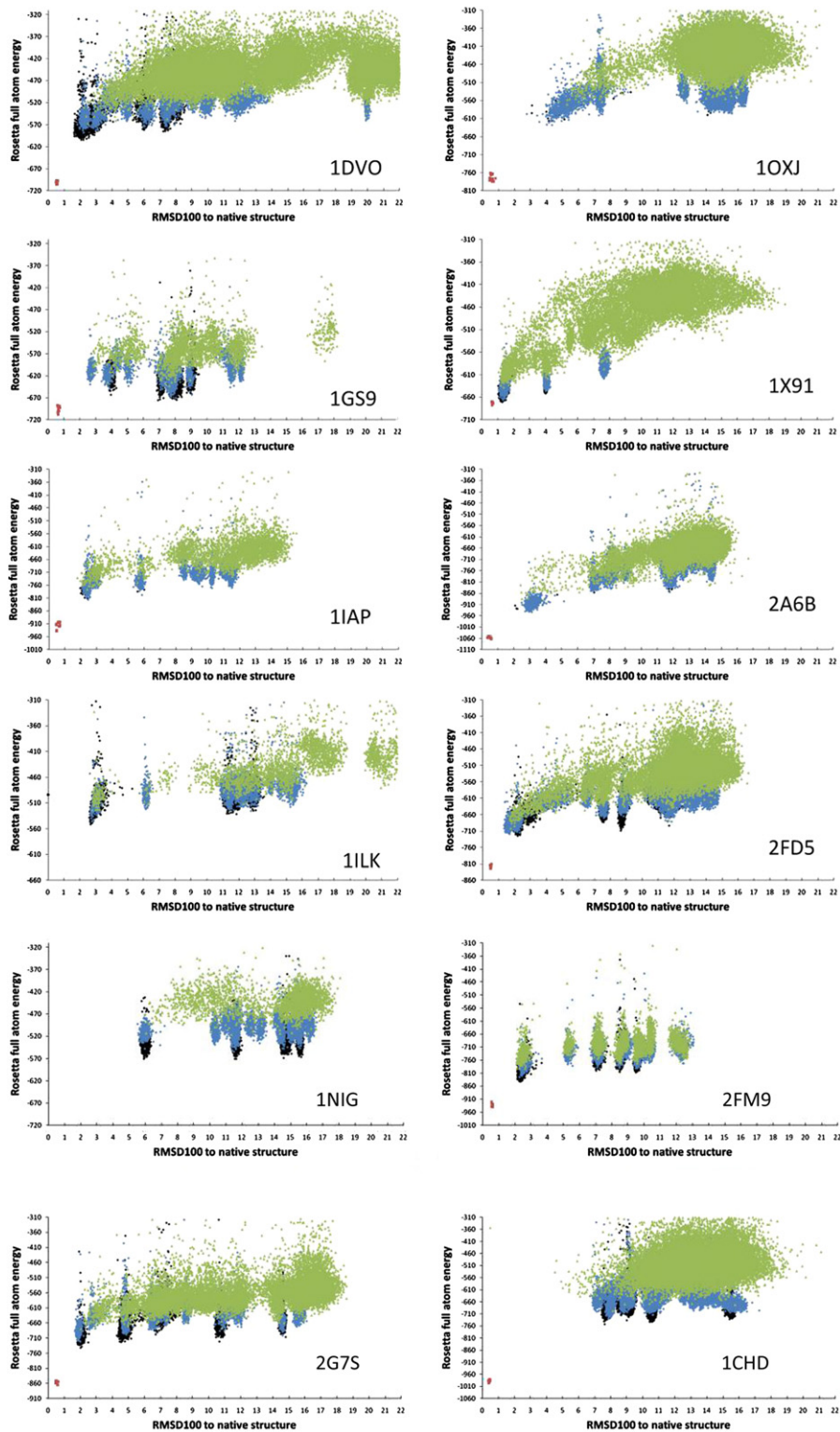
Despite the fact that the majority of the Rosetta-refined correct topologies actually scored best among all the refined models, there were still few cases in which some incorrect topologies scored better (see, for example, 1GS9, 1NIG, 2IGC, or 1JL1). A quality measure that could independently distinguish successful from unsuccessful cases would be desirable for situations when EM-Fold is used to build structures where the correct solution is not known. We developed a measure that is based on the depth of converged ensemble energy minimum (DoCEEM) presented in Raman et al. (2010). Instead of basing the measure on the mutual RMSD between models, DoCEEM is based on topology assignment. A topology is defined as a placement of specific SSEs into particular density rods noting the individual orientations of SSEs along the density rods' main axis (parallel or anti-parallel). Given the way the models are assembled and refined with EM-Fold, every generated model can be easily classified according to its topology. The DoCEEM is calculated as the energy difference between the median energy of the 10 lowest-scoring models with the same topology as the top-scoring model after the third round of Rosetta refinement and the median energy of the 10 lowest-scoring models with a topology different from the topology of the top-scoring model. Graphically speaking, the DoCEEM measures how deep the energy funnel is that separates the top-scoring model from all other models of different topologies. Ideally, the deeper the funnel, the more likely that the top-scoring model is the correct topology. The DoCEEM values for all 19 proteins refined with Rosetta are shown in Table 2. Using a cutoff of 0 Rosetta energy units, where negative values indicate success, the DoCEEM was able to correctly identify whether the top-scoring model was the correct topology for all but one protein (1GS9). These results suggest that deep energy funnels very likely correspond to models close to the native structure. This allows the user to employ the DoCEEM as a measure of how likely the algorithm found the true topology as the top-scoring topology.

Noise in Simulated Density Maps Causes Slight Performance Decrease in Loop-Building Steps

The major difference between the benchmark described so far and a real-world application is that the simulated density maps

used in the benchmark did not contain noise. This may not have had such a profound influence on the EM-Fold assembly and refinement steps, as the models only contained residues in SSEs, which are commonly well-defined even in experimental maps. However, noise may have a profound impact on the loop building and refinement procedure in Rosetta that relies explicitly on density in loop regions of the proteins. To test the performance of EM-Fold and Rosetta when confronted with noisy maps, noise was added randomly to the simulated maps until a cross-correlation of 0.8 between noise-free and noise-containing map was achieved. The procedure was described in Woetzel et al. (2011), and 0.8 had been established as a realistic value best mimicking experimental density maps. The influence of the actual density map on the assembly step is minimal, so it was not repeated. This way, it was also guaranteed that the same starting models were used and the comparison between noise-free and noisy maps was legitimate. All the parameters used were identical to the benchmark with the noise-free simulated density maps. Table S3 summarizes the results of the benchmark with the noisy maps in much the same way Table 1 does for the noise-free maps. The top-scoring 150 topologies from the assembly step were refined in EM-Fold using the noisy maps. The results of the refinement step confirmed that noise did not have a profound influence on the model quality during this step. The average RMSD100s were virtually identical (2.98 Å versus 2.91 Å). The average rank of the top-scoring correct topology decreased to 40 (from 24 in the noise-free case). It is more difficult to correctly rank models based on density cross-correlation when noise is present. This, however, does not present a major problem as long as more topologies are carried over to the Rosetta refinement rounds. The top 75 topologies after EM-Fold refinement were chosen for three rounds of Rosetta refinement. As a proof of principle, the proteins 2FD5 and 1CHD were also refined in Rosetta, despite their correct topologies having ranks worse than 75. The three rounds of Rosetta refinement ranked all the correct topologies among the top 10 scoring topologies, with the majority of the correct topologies scoring best in Rosetta's force field. Over the course of the three rounds, the RMSD100 values improved slightly on average but were about 0.7 Å worse than in the noise-free benchmark. This was not unexpected, as noise has the highest impact on loop regions. The overall quality of the models (10 models had an RMSD100 of <3 Å after the third round of Rosetta refinement) was still very good, albeit somewhat lower than in the noise-free benchmark. For illustration purposes, Figure 5 shows two examples of successful model building (2G7S and 1OZ9) as well as one protein (1NIG) for which the overall RMSD100 was well beyond the target of 3 Å. It can be concluded that the proposed combination of the new version of EM-Fold with Rosetta can be successfully applied to maps containing noise.

has 150 residues. The model shown has an RMSD100 of 1.36 Å over the full length of the protein and 0.99 Å over the residues in SSEs. (C) 1DVO has 152 residues. The model shown has an RMSD100 of 1.83 Å over the full length of the protein and 1.18 Å over the residues in SSEs. (D) Some models score well but exhibit a relatively high RMSD100. One example is the model for 1DVO, which is shown here and is represented by a yellow dot in (C). Overall, the agreement of the model with the native structure is good (left image). Closer evaluation (center image) reveals that one helix has been placed into the density rod in the wrong orientation (green). This leads to a high overall RMSD100 (5.94 Å) but preserves many of the contacts between correctly placed SSEs, thus ensuring a relatively good score. The image on the right shows a close-up view of that particular helix in the best-scoring model (represented by a red dot in C).



Web 3C

Figure 3. (Continued)

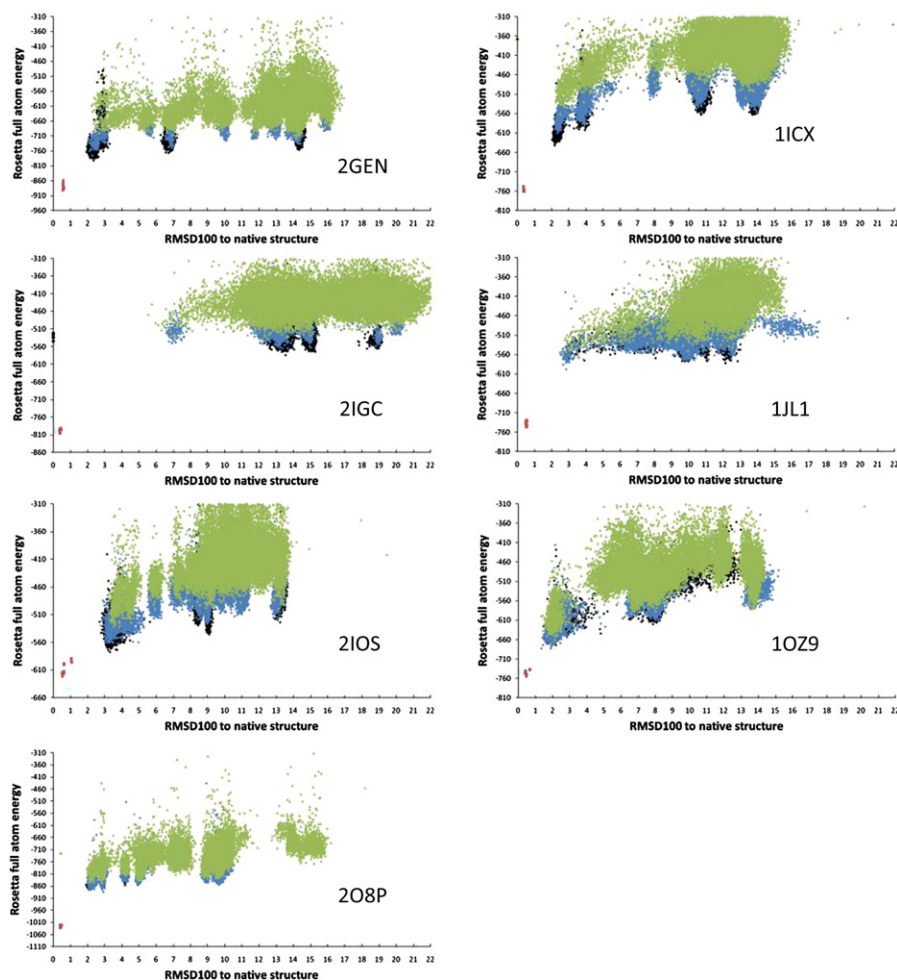


Figure 3. Gallery of Score versus RMSD100 Plots for All 19 Benchmark Proteins

Models from round 1 (green), round 2 (blue), and round 3 (black) of the Rosetta refinement are plotted. The native structure relaxed in Rosetta's force field is shown in violet for comparison. RMSD100s are calculated over all backbone atoms. The vast majority of the plots exhibit a clear funnel shape (i.e., models with low RMSD100 generally have lower scores than models with high RMSD100 values and vice versa). This feature is desirable in computational protein structure prediction because it makes model identification based on score possible.

See also [Figure S1](#).

EM-Fold and Rosetta Refinement Results in Models that Display Atomic Detail beyond that Present in the Density Map

In 14 out of 19 cases, Rosetta scored the correct topology as the best topology and resulted in a model with an RMSD100 below 2.8 Å. There was a fifteenth case (2IOS) where Rosetta scored the correct topology as the best topology, but the RMSD100 was somewhat worse. The correct topology was among the top four scoring topologies in all 19 proteins refined with Rosetta. The majority of the RMSD100 versus score plots showed a clear funnel-shape, indicating that Rosetta score correlated with model quality (see [Figure 3](#)). The vast majority of the best-scoring models is in excellent agreement with the native structure, as can be seen in [Figures 4](#) and [S2](#). For 70% of the successful benchmark cases, the best-scoring model had an RMSD100 below 3.0 Å, indicating that these models are accurate at atomic detail. EM-Fold identified the correct topology of a protein from

a medium-resolution density map in 70% of the cases. Refinement with Rosetta yielded models that were accurate at atomic detail in 70% of cases where the correct topology was identified. SSEs that were placed in the same density rod in at least 70% of the top 2,000 scoring models after the third round of Rosetta refinement were correctly placed with 97% confidence. Statistics over rotamer recovery revealed that after three rounds of Rosetta refinement, the best-scoring models recovered between 48% and 100% of the native rotamers in the protein core. The average rotamer recovery was 59%. In summary, many of the final refined models had a significant fraction of their native side chain conformations recovered correctly (see [Table 1](#)). This recovery was not based on information of side chain conformations in the medium-resolution density maps; rather, it was based on Rosetta's ability to correctly place side chains once the backbone conformation has approached the native structure. Hence the main achievement of the protocol

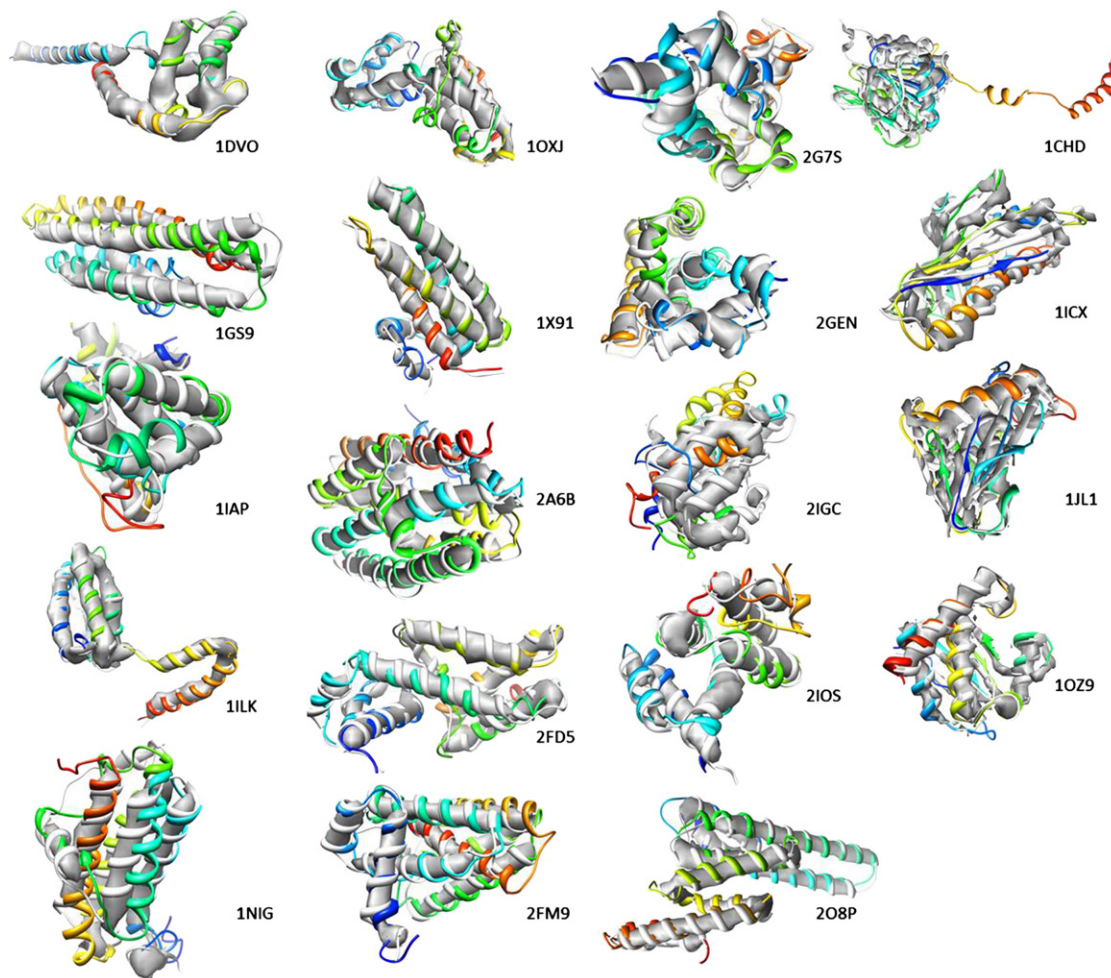


Figure 4. Gallery of Best-Scoring Models after Three Rounds of Rosetta Refinement for All 19 Benchmark Proteins

Superimposition of the models (rainbow-colored) with the original PDB structures (gray) and the simulated density maps are shown. Density maps were simulated at 5 Å resolution (β strand containing proteins) and at 7 Å resolution (α -helical proteins), respectively. The average RMSD100 of the models shown to the native structure is 2.96 Å, with RMSD100 values ranging from 1.29 to 6.93 Å.

See also Figure S2.

is that we were able to de novo build backbone models of the benchmark proteins guided by the density maps that are of sufficient quality to recover side chain conformations. The main reason for failure of assembly of the correct topology was inaccurate prediction of secondary structure. Primary obstacles for reaching atomic detail accuracy were a systematic sequence shift in SSEs that resulted in suboptimal starting models for Rosetta and long loop regions that were difficult to construct at atomic detail.

Applying the Protocol to Experimental Maps Yields Low RMSD Structures for SSE Model Parts and Atomic Detail in Favorable Cases

Due to the still limited number of experimental cryoEM density maps at resolutions between 5–7 Å of proteins for which high-resolution crystal structures also exist, the benchmark was performed on simulated density maps with added noise. To

test performance on experimental density maps, five proteins for which experimental maps and high-resolution structures are available were selected. Models were built into the bovine metarhodopsin cryoEM density map (Electron Microscopy Data Bank [EMDB] 1079 [Ruprecht et al., 2004], 5.5 Å resolution) and compared to the crystal structure of bovine rhodopsin (Protein Data Bank [PDB] ID 1GZM [Li et al., 2004]). Additionally, a model for proteins PrgH and PrgK was built into the subnanometer resolution structure from *Salmonella*'s needle complex (EMDB 1874 [Schraidt and Marlovits, 2011], subnanometer resolution) and compared to docked crystal structures of these components (PDB ID 2Y9J [Schraidt and Marlovits, 2011]). Finally, models for the 30S ribosomal proteins S15 and S20 were built into the ribosome cryoEM density map (EMDB 1829 [Bhushan et al., 2011], 5.6 Å resolution) and compared to the crystal structures of these proteins (PDB IDs 2WWLO and 2WWLT [Seidelt et al., 2009]).

Table 2. DoCEEM Analysis of Folding Results for All 19 Proteins Refined with Rosetta

Protein (PDB ID)	DoCEEM, REU ^a	Rosetta Round 3 Rank ^b
1DVO	-10.581	1
1GS9	5.087	4
1IAP	5.192	1
1ILK	-20.213	1
1NIG	1.843	2
1OXJ	-18.7	1
1X91	-6.57	1
2A6B	-73.26	1
2FD5	-24.344	1
2FM9	-32.335	1
2G7S	-11.522	1
2GEN	-23.86	1
2IGC	4.405	2
2IOS	-28.615	1
2O8P	-10.12	1
1CHD	5.764	4
1ICX	-52.623	1
1JL1	2.579	2
1OZ9	-49.456	1
1DVO	-10.581	1

^aDoCEEM values are calculated as the difference between the median energy of the 10 lowest-scoring models with the same topology as the top-scoring model after the third round of Rosetta refinement and the median energy of the 10 lowest-scoring models with a topology different from the topology of the top-scoring model. Negative values indicate that the mean energy of the top models of the top-scoring topology was lower than the mean energy of the top models of topologies different from the top-scoring topology.

^bThe DoCEEM values correlate with the Rosetta rank. This is important because in a nonbenchmark application, only DoCEEM values are available. REU, Rosetta energy units.

The same EM-Fold folding protocol applied to the benchmark proteins with simulated maps was used for the proteins with experimental density maps. While the ribosome proteins underwent three rounds of Rosetta refinement, only a single round of Rosetta refinement was performed for the other proteins in order to limit the computational resources needed. The results are summarized in Table 3. Additionally, Figure 6 and Figure S3 show the models after the EM-Fold assembly step, EM-Fold refinement step, and Rosetta loop building and refinement in context with the native structures and the experimental density maps. The correct topology scored within the top 40 topologies for all five proteins after the EM-Fold assembly step. For two of the proteins, the ranking of the correct topology improved slightly after EM-Fold refinement. The average RMSD100 of the correct topology models after the initial EM-Fold assembly step was 3.25 Å and improved to 2.86 Å after EM-Fold refinement. These values are in the same range as the RMSD100 values of the benchmark proteins using simulated density maps. Along with the results of the noisy maps benchmark, this is confirmation that EM-Fold also works well for experimental density maps. The results for the ribosomal proteins

S15 and S20 showed improvement over all rounds of model building. The final models exhibited RMSD100 values below 3 Å, and side chain conformations within the protein core were recovered especially for 2WWLO (see Figure 6G). The good results for these proteins are speculated to be mainly due to their high secondary structure content and the superb quality of the density map. For rhodopsin and salmonella proteins, the quality of the models after Rosetta loop building and refinement was lower than the model quality in the benchmarks with simulated density maps. The average RMSD100 over the full length of these proteins was 4.92 Å, while the average RMSD100 over residues in SSEs was 2.79 Å. The average deviation for these proteins was higher mainly because of long, floppy loop regions that were difficult to predict and because of a high amount of noise in loop regions in experimental density maps. While the RMSD100s were slightly higher than in the previously discussed benchmarks, even for these challenging cases EM-Fold proved to be a highly valuable tool to determine the correct topology of a protein based on the density map and predicted good models for protein structures even for experimental maps.

Conclusions

In summary, the combination of EM-Fold and Rosetta is a powerful tool for de novo folding of proteins into medium-resolution density maps. This report demonstrates that computational methods are capable of extending the information available from cryoEM density maps. We further demonstrate that the combination of EM-Fold and Rosetta can build an atomic model from a medium-resolution density map and the protein sequence. This will give researchers the opportunity to utilize medium-resolution density maps more effectively. This work also demonstrates that medium-resolution density maps can contribute valuable information regarding the true atomic resolution structure.

The results show that EM-Fold is the method of choice in cases when a medium-resolution density map is determined, SSEs are identifiable as density rods (either manually or using automated software), loop connectivity information is elusive, and no backbone trace or template start model are available. For higher-resolution maps that contain information on the SSE connectivity backbone, tracing techniques such as GORGON may be better suited (Baker et al., 2011). For maps at lower resolution, where SSEs are not clearly visible, techniques such as fitting of comparative models should be used (Topf et al., 2005).

It appears that all β strand proteins are more difficult to model accurately with EM-Fold. We attribute this observation to a combination of several effects: (1) accuracy of secondary structure prediction is reduced for β strands compared to α helices; (2) it is generally more difficult to de novo fold all β strand proteins due to the increased number of nonlocal contacts that have to be sampled (Bonneau et al., 2002); (3) all β strand proteins have a higher number of SSEs per residue, leading to a larger number of possible topologies that need to be sampled; and (4) all β strand proteins have a larger fraction of residues in loop regions. Because loop regions are modeled less accurately in the EM-Fold protocol, RMSD values for all β strand proteins are higher. Table 4 summarizes indicators a–d and relates

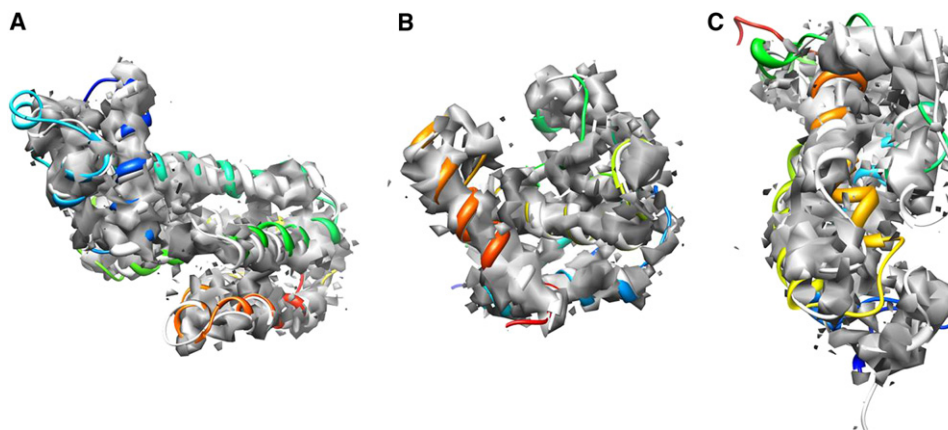


Figure 5. Three Protein Models after the Third Round of Rosetta Refinement in the Noisy Maps Benchmark

(A–C) Superimposition of the final models (rainbow-colored) of 2G7S (A), 1OZ9 (B), and 1NIG (C) with the original PDB structures (gray) as well as the noisy density maps are shown. Density maps were simulated at 5 Å resolution (β strand containing proteins) and at 7 Å resolution (α -helical proteins), respectively. (A) 2G7S has 194 residues. The model shown has an RMSD100 of 2.23 Å over the full length of the protein and 1.67 Å over the helical residues. (B) 1OZ9 has 150 residues. The model shown has an RMSD100 of 1.76 Å over the full length of the protein and 1.02 Å over the residues in SSEs. (C) 1NIG has 152 residues. The model shown has an RMSD100 of 7.31 Å over the full length of the protein and 4.97 Å over the residues in SSEs.

them to success in the benchmark. Although there is no single indicator of success, it seems that particularly the contact order and the fraction of residues in loop regions correlate with success. The average contact order in the benchmark set was about 13, while the average contact order for failures was 16. Similarly, the average fraction of residues in loop regions was 0.33 over the benchmark set and 0.42 for the failures. These quantities may give researchers using EM-Fold an indication of what the expected success may be.

In summary, substantial progress has been made since the initial EM-Fold release (Lindert et al., 2009). Because of its added features, the new version of EM-Fold can refine protein models to atomic detail accuracy in favorable cases, it is more tolerant to errors in secondary structure prediction, and can assemble

proteins that contain β strands. The consistent ability of predicting protein structure de novo and at atomic detail accuracy based on medium-resolution density maps is genuine progress in the field of cryoEM modeling techniques.

EXPERIMENTAL PROCEDURES

Folding Protocol

The folding protocol employed in this work is summarized in Figure 1. This basic protocol is based on the initial EM-Fold publication (Lindert et al., 2009). Several improvements have been added to this new version of EM-Fold. These will be discussed in greater detail here. Starting from the primary sequence of the protein, α helices and β strands are predicted using jufo (Meiler and Baker, 2003a; Meiler et al., 2001), PsiPred (Jones, 1999), and PROFphd (Rost and Sander, 1993a; Rost and Sander, 1993b; Rost and Sander, 1994). The predictions as well as their consensus are stored in a pool of SSEs (Figure 1A). The assembly step (Figure 1B) places SSEs from the pool into the density rods. It is assumed that the density for α helices and β strands is sufficiently different to exclusively place the correct SSEs into the individual density rods. In addition to the moves described by Lindert et al. (2009), growing and shrinking of SSEs is performed, helping alleviate some of the problems caused by incorrect secondary structure prediction. Lindert et al. (2009) showed that the secondary structure prediction algorithms generally underpredict the length of α helices. In the original implementation of EM-Fold, this was addressed by additional extended copies of α helices to the pool. Only the predicted SSEs are added to the pool in the new version. Subsequently, during the assembly Monte Carlo steps, SSEs are randomly grown or shrunk by up to two residues per step. The resizing is accompanied by a score that evaluates the agreement of secondary structure in the model with the predicted secondary structure. This ensures that SSEs remain in overall agreement with the predicted regions. Growth and shrinkage of SSEs has the potential to compensate for incorrect secondary structure prediction in a more dynamic way than the SSE pool used before. Models built in the assembly step are clustered. The best scoring clusters transition into the refinement step. The refinement step (Figure 1C) applies small translational and rotational perturbations to the SSEs in the model. When SSEs are placed into the density rods, they are idealized (i.e., perfectly straight). Some density rods, however, show at least a slight curvature. A new move that bends SSEs has been added to the refinement step. The center and amount of the bending are determined randomly. With bending in place, it is necessary to evaluate the agreement of the model with the density map. In the new implementation, this is done using a density cross correlation score. Scores

Table 3. Results of Benchmark on Experimental Density Maps from Rhodopsin, Salmonella, and Ribosome

Protein (PDB ID)	Size ^a	Rank/RMSD100, Å ^b		Rank/RMSD100 (RMSD100 SSEs), Å ^c
		EM-Fold Assembly	EM-Fold Refinement	
1GZM	349, 8, 0	37/2.26	35/2.79	2/4.60 (2.74)
2Y9JO	186, 4, 6	26/2.81	23/2.50	2/5.51 (2.71)
2Y9JZ	170, 4, 6	17/3.11	48/2.86	2/4.64 (2.91)
2WWLO	88, 4, 0	2/3.98	33/3.23	1/2.96 (2.01)
2WWLT	85, 3, 0	1/4.10	8/2.94	8/2.84 (2.41)

All RMSD100 values were determined over the backbone atoms N, C α , C, and O.

^aNumber of amino acids, number of α helices with at least 12 residues, and number of β strands with at least five residues.

^bThe rank of the correct topology model within all scored models as well as the RMSD100 of the correct topology model are given.

^cResults of Rosetta loop building and refinement using the experimental maps. The rank of the correct topology model within all scored models, the RMSD100 of the correct topology model, and the RMSD100 of the correct topology model over residues in SSEs (numbers in parentheses) are shown.

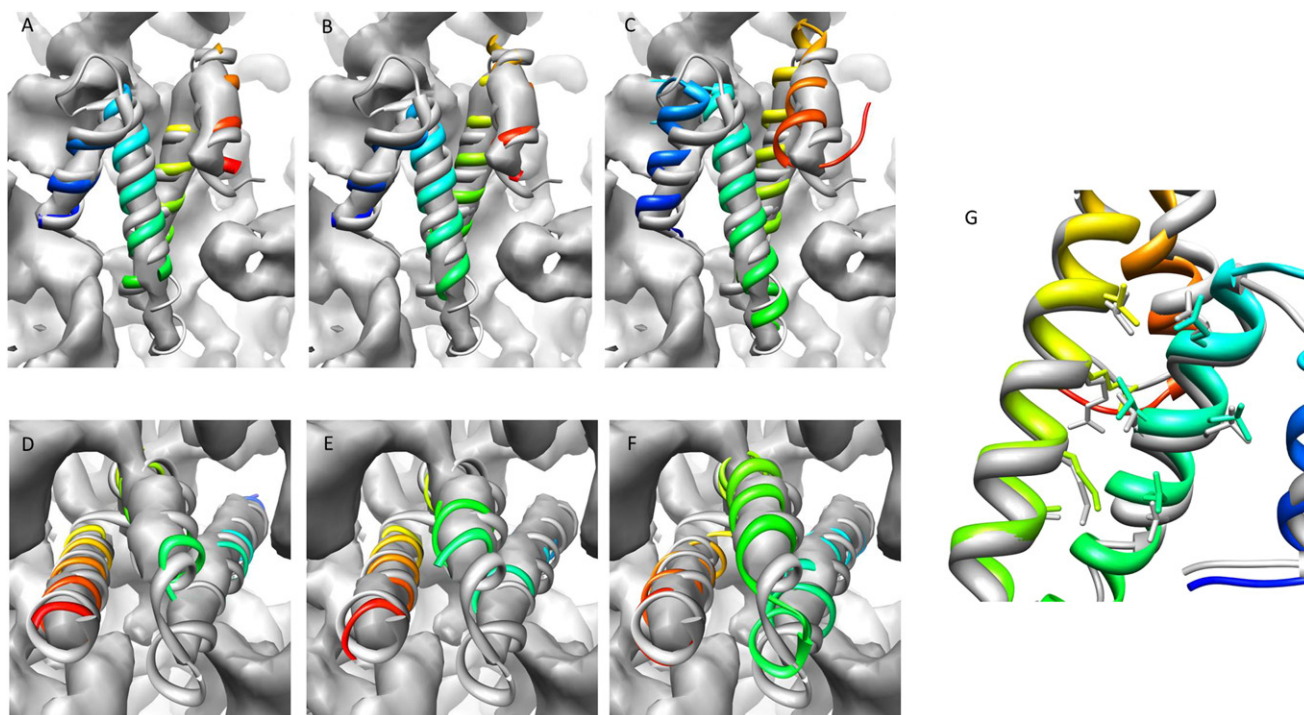


Figure 6. Results of Ribosome Benchmark Proteins with Experimental Density Maps

(A–F) Protein models after EM-Fold assembly step (A and D), after EM-Fold refinement step (B and E) and third round of Rosetta loop building and refinement (C and F) are shown for the two benchmark proteins. Superimposition of the models (rainbow-colored) of 2WWLO (A–C) and 2WWLT (D–F) with the original PDB structures (gray) as well as the experimental density maps are shown.

(G) Atomic detail recovered for 2WWLO after the third round of Rosetta refinement is shown.

See also Figure S3.

that evaluate solvation-free energy and residue-residue pairwise interaction within the protein are used in addition. A small number of top-scoring topologies identified in the refinement step will be used for loop building and refinement within Rosetta. Three rounds of Rosetta refinement (Figures 1D and 1E) build missing coordinates and refine the model further. The first round of Rosetta refinement (Figure 1D) builds missing loop coordinates guided by the density map. The executable used is `loopmodel.linuxgccrelease`. In the following two rounds (Figure 1E), regions of the models that agree least with the density map are identified (`loops_from_density.linuxgccrelease`) and rebuilt (`loopmodel.linuxgccrelease`). Each round includes a relaxation of the overall structure. An increasingly smaller number of topologies enter each of the three rounds of Rosetta refinement. After each round, the built models are clustered according to their topology and only the best scoring representative of the topologies will advance into the next round of refinement. After the third round of Rosetta refinement, the best scoring model is identified as the model for the protein structure.

Benchmark

A benchmark set of 20 α -helical and seven β sheet proteins was compiled from the PDB. The proteins ranged in size from 150 to 249 residues. Other selection criteria were secondary structure content of least 60% and the availability of high-resolution crystal structures for model comparison. The α -helical proteins contained between four and nine α helices, while the β sheet proteins contained between four and 10 β strands. Density maps at 7 and 5 Å resolution were simulated for the α -helical and β sheet proteins, respectively, using PDB2VOL from the SITUS package (Wriggers and Birmanns, 2001). A voxel spacing of 1.5 Å and Gaussian flattening was used. The rationale for this was that density maps at these resolutions will exhibit sufficient detail to identify density rods for α helices and β strands, respectively. Additionally, density maps at 7 and 5 Å resolution that had noise added were simulated

using the PDB2DENSITY application of the BCL. Noise was added randomly until a cross-correlation between noisy and noise-free maps dropped below 0.8. Jufo, PsiPred, and PROFphd were used to predict SSEs for the consensus pool. A three-state model (helix, strand, coil) was used for the pool. α Helices with 12 or more predicted residues and β strands with at least five predicted residues were added to the pool. Shorter secondary elements were omitted from the initial EM-Fold assembly and added in the later Rosetta refinement step. In the assembly step, 50,000 models (2,000 rejected Monte Carlo steps) were built for each of the 27 proteins. Building one model took approximately 60 s on a single 2.4 GHz Quad-Core AMD Opteron Processor. Building 50,000 models took approximately 2–3 hr on a 400 core cluster. The 50,000 models were clustered into topologies according to their placement of particular stretches of sequence into density rods. The topologies were ranked by the overall score and the top-scoring models within each of the top-scoring 150 topologies advance to the refinement step. If the correct topology was not identified within the top 150 scoring topologies, this protein counted as a failure in the benchmark. Quality of the models is determined by calculating the RMSD100 (Carugo and Pongor, 2001) values over the backbone atoms N, C α , C, and O. For each of the top 150 scoring topologies from the assembly step, 500 refined models were built in the refinement step. Again, building one model took approximately 60 s on a single 2.4 GHz Quad-Core AMD Opteron Processor. Building 75,000 refined models took approximately 3–4 hr on a 400 core cluster. These models were ranked by their refinement score, and the top-scoring 50 (75 in the noisy map benchmark) topologies after refinement step were identified. These 50 models served as input for the first round of Rosetta refinement. Rosetta refinement was performed on the ACCRE computer cluster (2.4 GHz Quad-Core AMD Opteron Processors). Timing of Rosetta refinement depended heavily on the size of the protein, the size of the loop regions, and the size of the density map. Generally, it can be assumed that a single round of refinement for one

Table 4. Summary of Protein Statistics and Benchmark Performance

Protein (PDB ID)	Residues	SSPred Accuracy	Contact Order	SSE per Residue	Fraction Loop	Success
1DVO	152	0.74	9.28	0.03	0.26	Atomic resolution
1GS9	165	0.87	11.67	0.02	0.3	Topology
1IAP	211	0.79	14.86	0.03	0.39	Atomic resolution
1ILK	151	0.81	6.08	0.03	0.26	Atomic resolution
1NIG	152	0.69	12.79	0.03	0.32	Topology
1OXJ	173	0.80	8.47	0.02	0.32	Topology
1X91	153	0.78	11.21	0.03	0.24	Atomic resolution
1Z3Y	238	0.78	12.12	0.03	0.42	Failure
2A6B	234	0.86	24.06	0.03	0.28	Atomic resolution
2FD5	180	0.80	11.18	0.03	0.26	Atomic resolution
2FM9	215	0.82	13	0.04	0.28	Atomic resolution
2FQ4	192	0.83	9.06	0.04	0.34	Failure
2G7S	194	0.76	11.08	0.03	0.22	Atomic resolution
2GEN	197	0.80	10.75	0.04	0.23	Atomic resolution
2IGC	164	0.69	11.29	0.02	0.31	Topology
2IOS	150	0.80	8.52	0.04	0.34	Topology
2IU1	208	0.77	13.03	0.02	0.38	Failure
2NR7	195	0.75	13.91	0.03	0.33	Failure
2O8P	227	0.79	8.53	0.04	0.19	Atomic resolution
2QK1	249	0.81	9.05	0.04	0.32	Failure
1BJ7	156	0.79	20.86	0.06	0.36	Failure
1CHD	203	0.73	19.04	0.04	0.43	Topology
1ICX	155	0.77	18.63	0.05	0.36	Atomic resolution
1JL1	155	0.77	11.35	0.05	0.32	Atomic resolution
1OZ9	150	0.80	7.1	0.06	0.32	Atomic resolution
1WBA	175	0.80	21.541	0.06	0.54	Failure
2QVK	192	0.72	28.262	0.04	0.69	Failure

For each of the proteins used in the benchmark, the number of residues, the accuracy of the consensus secondary structure prediction, the absolute contact order, the number of SSEs per residue, and the fraction of residues that are in loop regions are shown. Additionally, it is labeled whether the benchmark protocol managed to predict the protein's structure to atomic detail, whether atomic detail was not achieved but the protein topology has been predicted correctly, or whether EM-Fold failed to predict the structure of this protein.

protein takes about 24 hr on a 400 core cluster. In the first round of refinement, Rosetta built loop models for these 50 topologies guided by the density map. The models were clustered according to topology, and the top-scoring models from the top 15 clusters entered round 2 of Rosetta refinement. Using the `loops_from_density.linuxgccrelease` executable, regions within these 15 proteins that agreed least with the density map were identified. These regions were subsequently rebuilt using the guidance of the density map. The top-scoring five topologies after the second round of Rosetta refinement were used as starting models in the third round of refinement. After three rounds of refinement, the best scoring model was evaluated.

Software Availability

EM-Fold is part of the BCL software library developed in the Meiler laboratory. It is supported for Linux, Windows, and Mac environments. EM-Fold is freely available to the scientific community at <http://bclcommons.vueinnovations.com/licensing>.

SUPPLEMENTAL INFORMATION

Supplemental Information includes three figures and three tables and can be found with this article online at [doi:10.1016/j.str.2012.01.023](https://doi.org/10.1016/j.str.2012.01.023).

ACKNOWLEDGMENTS

This research was supported by grants to J.M. (NSF 0742762 "CAREER: Cryo-EM guided de novo Protein Fold Elucidation" and NIH 1R01GM080403 "Membrane Protein Structure Elucidation from sparse NMR data" [KAMP]).

Received: May 5, 2011

Revised: January 23, 2012

Accepted: January 26, 2012

Published: March 6, 2012

REFERENCES

- Alexander, N., Bortolus, M., Al-Mestarihi, A., Mchourab, H., and Meiler, J. (2008). De novo high-resolution protein structure determination from sparse spin-labeling EPR data. *Structure* 16, 181–195.
- Baker, M.L., Ju, T., and Chiu, W. (2007). Identification of secondary structure elements in intermediate-resolution density maps. *Structure* 15, 7–19.
- Baker, M.L., Abeysinghe, S.S., Schuh, S., Coleman, R.A., Abrams, A., Marsh, M.P., Hryc, C.F., Ruths, T., Chiu, W., and Ju, T. (2011). Modeling protein structure at near atomic resolutions with Gorgon. *J. Struct. Biol.* 174, 360–373.

- Bhushan, S., Hoffmann, T., Seidelt, B., Frauenfeld, J., Mielke, T., Berninghausen, O., Wilson, D.N., and Beckmann, R. (2011). SecM-stalled ribosomes adopt an altered geometry at the peptidyl transferase center. *PLoS Biol.* 9, e1000581.
- Bonneau, R., Ruczinski, I., Tsai, J., and Baker, D. (2002). Contact order and ab initio protein structure prediction. *Protein Sci.* 11, 1937–1944.
- Bowers, P.M., Strauss, C.E., and Baker, D. (2000). De novo protein structure determination using sparse NMR data. *J. Biomol. NMR* 18, 311–318.
- Bradley, P., Misura, K.M., and Baker, D. (2005). Toward high-resolution de novo structure prediction for small proteins. *Science* 309, 1868–1871.
- Carugo, O., and Pongor, S. (2001). A normalized root-mean-square distance for comparing protein three-dimensional structures. *Protein Sci.* 10, 1470–1473.
- Cong, Y., Baker, M.L., Jakana, J., Woolford, D., Miller, E.J., Reissmann, S., Kumar, R.N., Redding-Johanson, A.M., Bath, T.S., Mukhopadhyay, A., et al. (2010). 4.0-Å resolution cryo-EM structure of the mammalian chaperonin TRiC/CCT reveals its unique subunit arrangement. *Proc. Natl. Acad. Sci. USA* 107, 4967–4972.
- Dal Palù, A., He, J., Pontelli, E., and Lu, Y. (2006). Identification of alpha-helices from low resolution protein density maps. *Comput. Syst. Bioinformatics Conf.* 5, 89–98.
- DiMaio, F., Tyka, M.D., Baker, M.L., Chiu, W., and Baker, D. (2009). Refinement of protein structures into low-resolution density maps using rosetta. *J. Mol. Biol.* 392, 181–190.
- Hanson, S.M., Dawson, E.S., Francis, D.J., Van Eps, N., Klug, C.S., Hubbell, W.L., Meiler, J., and Gurevich, V.V. (2008). A model for the solution structure of the rod arrestin tetramer. *Structure* 16, 924–934.
- Hirst, S.J., Alexander, N., McHaourab, H.S., and Meiler, J. (2011). RosettaEPR: an integrated tool for protein structure determination from sparse EPR data. *J. Struct. Biol.* 173, 506–514.
- Jiang, W., Baker, M.L., Ludtke, S.J., and Chiu, W. (2001). Bridging the information gap: computational tools for intermediate resolution structure interpretation. *J. Mol. Biol.* 308, 1033–1044.
- Jones, D.T. (1999). Protein secondary structure prediction based on position-specific scoring matrices. *J. Mol. Biol.* 292, 195–202.
- Kong, Y., Zhang, X., Baker, T.S., and Ma, J. (2004). A structural-informatics approach for tracing beta-sheets: building pseudo-C(alpha) traces for beta-strands in intermediate-resolution density maps. *J. Mol. Biol.* 339, 117–130.
- Li, J., Edwards, P.C., Burghammer, M., Villa, C., and Schertler, G.F. (2004). Structure of bovine rhodopsin in a trigonal crystal form. *J. Mol. Biol.* 343, 1409–1438.
- Lindert, S., Staritzbichler, R., Wötzel, N., Karakaş, M., Stewart, P.L., and Meiler, J. (2009). EM-fold: de novo folding of α -helical proteins guided by intermediate-resolution electron microscopy density maps. *Structure* 17, 990–1003.
- Liu, H., Jin, L., Koh, S.B., Atanasov, I., Schein, S., Wu, L., and Zhou, Z.H. (2010a). Atomic structure of human adenovirus by cryo-EM reveals interactions among protein networks. *Science* 329, 1038–1043.
- Liu, X., Zhang, Q., Murata, K., Baker, M.L., Sullivan, M.B., Fu, C., Dougherty, M.T., Schmid, M.F., Osburne, M.S., Chisholm, S.W., and Chiu, W. (2010b). Structural changes in a marine podovirus associated with release of its genome into *Prochlorococcus*. *Nat. Struct. Mol. Biol.* 17, 830–836.
- Ludtke, S.J., Chen, D.-H., Song, J.-L., Chuang, D.T., and Chiu, W. (2004). Seeing GroEL at 6 Å resolution by single particle electron cryomicroscopy. *Structure* 12, 1129–1136.
- Ludtke, S.J., Baker, M.L., Chen, D.H., Song, J.L., Chuang, D.T., and Chiu, W. (2008). De novo backbone trace of GroEL from single particle electron cryomicroscopy. *Structure* 16, 441–448.
- Meiler, J., and Baker, D. (2003a). Coupled prediction of protein secondary and tertiary structure. *Proc. Natl. Acad. Sci. USA* 100, 12105–12110.
- Meiler, J., and Baker, D. (2003b). Rapid protein fold determination using unassigned NMR data. *Proc. Natl. Acad. Sci. USA* 100, 15404–15409.
- Meiler, J., and Baker, D. (2005). The fumarate sensor DcuS: progress in rapid protein fold elucidation by combining protein structure prediction methods with NMR spectroscopy. *J. Magn. Reson.* 173, 310–316.
- Meiler, J., Müller, M., Zeidler, A., and Schmäscke, F. (2001). Generation and evaluation of dimension-reduced amino acid parameter representations by artificial neural networks. *J. Mol. Model.* 7, 360–369.
- Min, G.W., Wang, H.B., Sun, T.T., and Kong, X.P. (2006). Structural basis for tetraspanin functions as revealed by the cryo-EM structure of uroplakin complexes at 6-Å resolution. *J. Cell Biol.* 173, 975–983.
- Raman, S., Lange, O.F., Rossi, P., Tyka, M., Wang, X., Aramini, J., Liu, G., Ramelot, T.A., Eletsky, A., Szyperski, T., et al. (2010). NMR structure determination for larger proteins using backbone-only data. *Science* 327, 1014–1018.
- Rohl, C.A., and Baker, D. (2002). De novo determination of protein backbone structure from residual dipolar couplings using Rosetta. *J. Am. Chem. Soc.* 124, 2723–2729.
- Rost, B., and Sander, C. (1993a). Improved prediction of protein secondary structure by use of sequence profiles and neural networks. *Proc. Natl. Acad. Sci. USA* 90, 7558–7562.
- Rost, B., and Sander, C. (1993b). Prediction of protein secondary structure at better than 70% accuracy. *J. Mol. Biol.* 232, 584–599.
- Rost, B., and Sander, C. (1994). Combining evolutionary information and neural networks to predict protein secondary structure. *Proteins* 19, 55–72.
- Ruprecht, J.J., Mielke, T., Vogel, R., Villa, C., and Schertler, G.F. (2004). Electron crystallography reveals the structure of metarhodopsin I. *EMBO J.* 23, 3609–3620.
- Saban, S.D., Silvestry, M., Nemerow, G.R., and Stewart, P.L. (2006). Visualization of alpha-helices in a 6-angstrom resolution cryoelectron microscopy structure of adenovirus allows refinement of capsid protein assignments. *J. Virol.* 80, 12049–12059.
- Schraidt, O., and Marlovits, T.C. (2011). Three-dimensional model of Salmonella's needle complex at subnanometer resolution. *Science* 331, 1192–1195.
- Seidelt, B., Innis, C.A., Wilson, D.N., Gartmann, M., Armache, J.P., Villa, E., Trabuco, L.G., Becker, T., Mielke, T., Schulten, K., et al. (2009). Structural insight into nascent polypeptide chain-mediated translational stalling. *Science* 326, 1412–1415.
- Serysheva, I.I., Ludtke, S.J., Baker, M.L., Cong, Y., Topf, M., Eramian, D., Sali, A., Hamilton, S.L., and Chiu, W. (2008). Subnanometer-resolution electron cryomicroscopy-based domain models for the cytoplasmic region of skeletal muscle RyR channel. *Proc. Natl. Acad. Sci. USA* 105, 9610–9615.
- Sibanda, B.L., Chirgadze, D.Y., and Blundell, T.L. (2010). Crystal structure of DNA-PKcs reveals a large open-ring cradle comprised of HEAT repeats. *Nature* 463, 118–121.
- Topf, M., Baker, M.L., John, B., Chiu, W., and Sali, A. (2005). Structural characterization of components of protein assemblies by comparative modeling and electron cryo-microscopy. *J. Struct. Biol.* 149, 191–203.
- Topf, M., Baker, M.L., Marti-Renom, M.A., Chiu, W., and Sali, A. (2006). Refinement of protein structures by iterative comparative modeling and CryoEM density fitting. *J. Mol. Biol.* 357, 1655–1668.
- Villa, E., Sengupta, J., Trabuco, L.G., LeBarron, J., Baxter, W.T., Shaikh, T.R., Grassucci, R.A., Nissen, P., Ehrenberg, M., Schulten, K., and Frank, J. (2009). Ribosome-induced changes in elongation factor Tu conformation control GTP hydrolysis. *Proc. Natl. Acad. Sci. USA* 106, 1063–1068.
- Ward, A., Reyes, C.L., Yu, J., Roth, C.B., and Chang, G. (2007). Flexibility in the ABC transporter MsbA: Alternating access with a twist. *Proc. Natl. Acad. Sci. USA* 104, 19005–19010.
- Wotzel, N., Lindert, S., Stewart, P.L., and Meiler, J. (2011). BCL::EM-Fit: Rigid body fitting of atomic structures into density maps using geometric hashing and real space refinement. *J. Struct. Biol.* 175, 264–276.

- Wriggers, W., and Birmanns, S. (2001). Using situs for flexible and rigid-body fitting of multiresolution single-molecule data. *J. Struct. Biol.* *133*, 193–202.
- Zhang, R., Hryc, C.F., Cong, Y., Liu, X., Jakana, J., Gorchakov, R., Baker, M.L., Weaver, S.C., and Chiu, W. (2011). 4.4 Å cryo-EM structure of an enveloped alphavirus Venezuelan equine encephalitis virus. *EMBO J.* *30*, 3854–3863.
- Zhou, Z.H. (2008). Towards atomic resolution structural determination by single-particle cryo-electron microscopy. *Curr. Opin. Struct. Biol.* *18*, 218–228.
- Zhu, J., Cheng, L., Fang, Q., Zhou, Z.H., and Honig, B. (2010). Building and refining protein models within cryo-electron microscopy density maps based on homology modeling and multiscale structure refinement. *J. Mol. Biol.* *397*, 835–851.