

## Macromolecular modeling and design in Rosetta: new methods and frameworks

Koehler Leman, Julia\* [1, 2], Weitzner, Brian D [3, 4, 5, 6], Lewis, Steven M [7, 8, 9],  
RosettaCommons consortium, Bonneau, Richard [1, 2, 60, 61]

### RosettaCommons consortium:

Adolf-Bryfogle, Jared [10], Alam, Nawsad [11], Alford, Rebecca F [3], Aprahamian, Melanie [12], Baker, David [4, 5], Barlow, Kyle A [13], Basanta, Benjamin [4, 14], Bender, Brian J [15], Blacklock, Kristin [16], Bonet, Jaume [17, 18], Boyken, Scott [5, 6], Bradley, Phil [19], Bystroff, Chris [20], Conway, Patrick [4], Cooper, Seth [21], Correia, Bruno E [17, 18], Coventry, Brian [4], Das, Rhiju [22], De Jong, René M [23], DiMaio, Frank [4, 5], Dsilva, Lorna [21], Dunbrack, Roland [24], Ford, Alex [4], Frenz, Brandon [9], Fu, Darwin Y [25], Geniesse, Caleb [22], Goldschmidt, Lukasz [4], Gowthaman, Ragul [26, 27], Gray, Jeffrey J [3], Gront, Dominik [28], Guffy, Sharon [7], Horowitz, Scott [29, 30], Huang, Po-Ssu [4], Huber, Thomas [31], Jacobs, Tim M [32], Jeliaskov, Jeliasko R [33], Johnson, David K [34], Kappel, Kalli [35], Karanicolas, John [24], Khakzad, Hamed [18, 36, 37], Khar, Karen R [34], Khare, Sagar D [16, 38, 52, 53, 54], Khatib, Firas [39], Khrumushin, Alisa [11], King, Indigo C [4, 9], Kleffner, Robert [21], Koepnick, Brian [4], Kortemme, Tanja [40], Kuenze, Georg [25, 41], Kuhlman, Brian [7], Kuroda, Daisuke [42, 43], Labonte, Jason W [3, 44], Lapidoth, Gideon [45], Leaver-Fay, Andrew [7], Lindert, Steffen [12], Linsky, Thomas [4, 5], London, Nir [11], Lubin, Joseph H [3], Lyskov, Sergey [3], Maguire, Jack [32], Malmström, Lars [18, 36, 37, 46], Marcos, Enrique [4, 47], Marcu, Orly [11], Marze, Nicholas A [3], Meiler, Jens [41, 48, 49], Moretti, Rocco [25], Mulligan, Vikram Khipple [1, 4, 5], Nerli, Santrupti [50], Norn, Christoffer [45], Ó'Conchúir, Shane [40], Ollikainen, Noah [40], Ovchinnikov, Sergey [4, 5, 51], Pacella, Michael S [3], Pan, Xingjie [40], Park, Hahnbeom [4], Pavlovicz, Ryan E [4, 5], Pethe, Manasi [52, 53], Pierce, Brian G [26, 27], Pilla, Kala Bharath [31], Raveh, Barak [11], Renfrew, P Douglas [1], Roy Burman, Shourya S [3], Rubenstein, Aliza [16, 54], Sauer, Marion F [55], Scheck, Andreas [17, 18], Schief, William [10], Schueler-Furman, Ora [11], Sedan, Yuval [11], Sevy, Alexander M [55], Sgourakis, Nikolaos G [56], Shi, Lei [4, 5], Siegel, Justin [57, 58, 59], Silva, Daniel-Adriano [4], Smith, Shannon [25], Song, Yifan [4, 5], Stein, Amelie [40], Szegedy, Maria [38], Teets, Frank D [7], Thyme, Summer B [4], Wang, Ray Yu-Ruei [4], Watkins, Andrew [22], Zimmerman, Lior [11]

1 Center for Computational Biology, Flatiron Institute, Simons Foundation, New York, NY 10010, USA

2 Dept of Biology, New York University, New York, 10003, NY, USA

3 Dept of Chemical and Biomolecular Engineering, Johns Hopkins University, Baltimore, MD 21218, USA

4 Dept of Biochemistry, University of Washington, Seattle, WA 98195, USA

5 Institute for Protein Design, University of Washington, Seattle, WA 98195, USA

6 Lyell Immunopharma Inc., Seattle, WA 98109

7 Dept of Biochemistry and Biophysics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA

8 Dept of Biochemistry, Duke University, Durham, NC 27710, USA

9 Cyrus Biotechnology, Seattle, WA 98101, USA

10 Dept of Immunology and Microbiology, The Scripps Research Institute, La Jolla, CA, USA

11 Dept of Microbiology and Molecular Genetics, IMRIC, Ein Kerem Faculty of Medicine, Hebrew University of Jerusalem, 91120, Jerusalem, Israel

12 Dept of Chemistry and Biochemistry, Ohio State University, Columbus, OH, 43210, USA

13 Graduate Program in Bioinformatics, University of California San Francisco, CA 94158, USA

14 Biological Physics Structure and Design PhD Program, University of Washington, Seattle, WA 98195, USA

15 Department of Pharmacology, Vanderbilt University, Nashville, TN 37232, USA

16 Institute of Quantitative Biomedicine, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA

17 Institute of Bioengineering, École Polytechnique Fédérale de Lausanne, Lausanne, Switzerland

18 Swiss Institute of Bioinformatics, Lausanne, Switzerland

19 Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA

20 Dept of Biological Sciences, Rensselaer Polytechnic Institute, Troy, NY, 12180, USA

21 Khoury College of Computer Sciences, Northeastern University, Boston, MA 02115, USA

- 22 Dept of Biochemistry, Stanford University School of Medicine, Stanford CA 94305, USA
- 23 DSM Biotechnology Center, 2613 AX Delft, The Netherlands
- 24 Institute for Cancer Research, Fox Chase Cancer Center, Philadelphia PA 19111, USA
- 25 Dept of Chemistry, Vanderbilt University, Nashville, TN 37232, USA
- 26 University of Maryland Institute for Bioscience and Biotechnology Research, Rockville, MD 20850, USA
- 27 Dept of Cell Biology and Molecular Genetics, University of Maryland, College Park, MD 20742, USA
- 28 Faculty of Chemistry, Biological and Chemical Research Centre, University of Warsaw  
Żwirki i Wigury 101, 02-089 Warsaw
- 29 Dept of Chemistry & Biochemistry, University of Denver, Denver, CO 80208, USA
- 30 The Knobel Institute for Healthy Aging, University of Denver, Denver, CO 80208, USA
- 31 Research School of Chemistry, Australian National University, Canberra ACT 2601, Australia
- 32 Program in Bioinformatics and Computational Biology, Dept of Biochemistry and Biophysics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, USA
- 33 Program in Molecular Biophysics, Johns Hopkins University, Baltimore, MD 21218, USA
- 34 Center for Computational Biology, University of Kansas, Lawrence, KS 66047, USA
- 35 Biophysics Program, Stanford University, Stanford CA 94305, USA
- 36 Institute for Computational Science, University of Zurich, CH-8057 Zurich, Switzerland
- 37 S3IT, University of Zurich, CH-8057 Zurich, Switzerland
- 38 Dept of Chemistry and Chemical Biology, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA
- 39 Dept of Computer and Information Science, University of Massachusetts Dartmouth, Dartmouth, MA 02747, USA
- 40 Dept of Bioengineering and Therapeutic Sciences, University of California San Francisco, CA 94158, USA
- 41 Center for Structural Biology, Vanderbilt University, Nashville, TN 37232, USA
- 42 Medical Device Development and Regulation Research Center, School of Engineering, University of Tokyo, Tokyo 113-8656, Japan
- 43 Dept of Bioengineering, School of Engineering, University of Tokyo, Tokyo 113-8656, Japan
- 44 Dept of Chemistry, Franklin & Marshall College, Lancaster, PA 17604, USA
- 45 Dept of Biomolecular Sciences, Weizmann Institute of Science, Rehovot, 76100, Israel
- 46 Division of Infection Medicine, Dept of Clinical Sciences Lund, Faculty of Medicine, Lund University, SE-22184, Lund, Sweden
- 47 Institute for Research in Biomedicine Barcelona, The Barcelona Institute of Science and Technology, 08028 Barcelona, Spain
- 48 Depts of Chemistry, Pharmacology and Biomedical Informatics, Vanderbilt University, Nashville, TN 37232, USA
- 49 Institute for Chemical Biology, Vanderbilt University, Nashville, TN 37232, USA
- 50 Dept of Computer Science, University of California Santa Cruz, CA, USA
- 51 Molecular and Cellular Biology Program, University of Washington, Seattle, WA 98195, USA
- 52 Dept of Chemistry and Chemical Biology, The State University of New Jersey, Piscataway, NJ 08854, USA
- 53 Center for Integrative Proteomics Research, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA
- 54 Computational Biology and Molecular Biophysics Program, Rutgers, The State University of New Jersey, Piscataway, NJ 08854, USA
- 55 Chemical and Physical Biology Program, Vanderbilt Vaccine Center, Vanderbilt University, Nashville, TN 37235, USA
- 56 Dept of Chemistry and Biochemistry, University of California Santa Cruz, CA, USA
- 57 Dept of Chemistry, University of California, Davis, CA, USA
- 58 Dept of Biochemistry and Molecular Medicine, University of California, Davis, CA, USA
- 59 Genome Center, University of California, Davis, CA, USA
- 60 Dept of Computer Science, New York University, New York, 10003, NY, USA
- 61 Center for Data Science, New York University, New York, 10003, NY, USA

\* Corresponding author: Julia Koehler Leman ([Julia.koehler.leman@nyu.edu](mailto:Julia.koehler.leman@nyu.edu))

**Abstract**

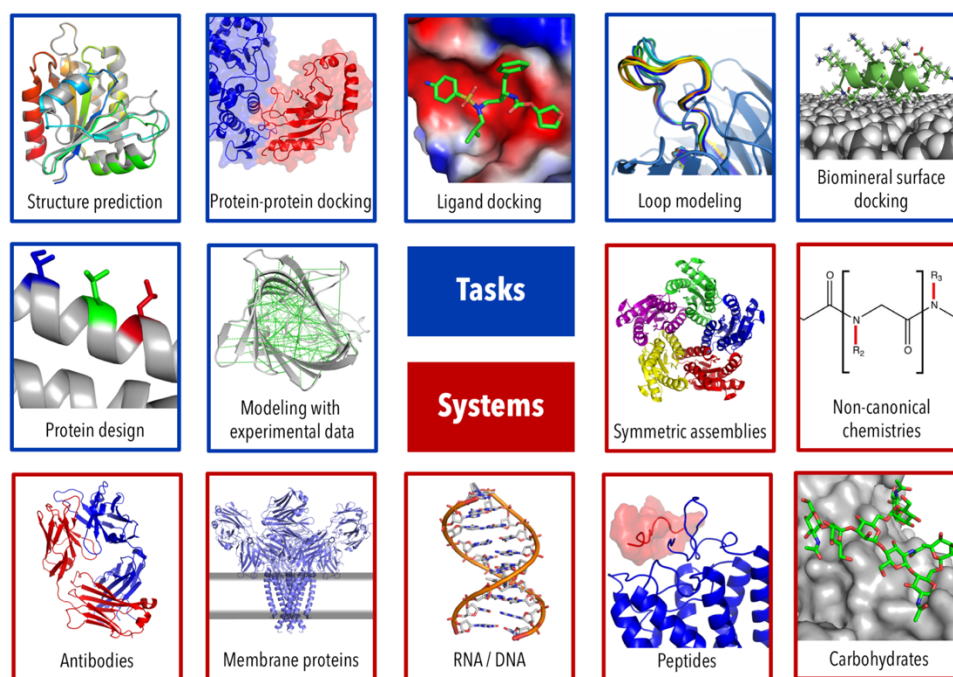
The Rosetta software suite for macromolecular modeling, docking, and design is widely used in pharmaceutical, industrial, academic, non-profit, and government laboratories. Considering its broad modeling capabilities, Rosetta consistently ranks highly when compared to other leading methods created for highly specialized protein modeling and design tasks. Developed for over two decades by a global community of scientists at more than 60 institutions, Rosetta has undergone multiple refactorings, and now comprises over three million lines of code. Here we discuss the methods developed in the last five years, involving the latest protocols for structure prediction, protein–protein and protein–small molecule docking, protein structure and interface design, loop modeling, the incorporation of various types of experimental data, and modeling of peptides, antibodies and other proteins in the immune system, nucleic acids, non-standard amino acids, carbohydrates, and membrane proteins. We briefly discuss improvements to the energy function, user interfaces, and usability of the software. Rosetta is available at [www.rosettacommons.org](http://www.rosettacommons.org).

## Introduction

Development of Rosetta started in the mid-1990s for protein structure prediction and to gain insights into the protein folding problem, which remains a grand challenge of structural biology. Over time, the number of applications grew to address a wider array of modeling tasks, ranging from protein–protein or –small molecule docking to incorporating NMR data, loop modeling, protein design, and interaction with peptides and nucleic acids (Figure 1). Over the 20 plus years since those early beginnings, the community of developers and scientists, the RosettaCommons, grew from a single academic laboratory to laboratories at over 60 institutions around the globe<sup>1</sup>. Rosetta has undergone several transitions, including in programming language and implementation, with the latest protocols based on Rosetta3, which was first released in 2008<sup>2</sup>. Rosetta's energy function has been continuously improved over its lifetime, detailed descriptions of which can be found in references <sup>3</sup> and <sup>4</sup>. As the Rosetta community grew, efforts to improve usability, interfaces to the code, and documentation have drastically improved usability and modular application to new problems. As part of our sustained focus on reproducibility and usability, we developed several interfaces, (PyRosetta<sup>5</sup>, RosettaScripts<sup>6</sup>, Foldit<sup>7</sup>) and emphasized publishing protocol captures<sup>8</sup> that accompany manuscripts to improve accessibility, user friendliness and scientific reproducibility. As the software's interfaces have grown more versatile, development has accelerated and branched in many directions. However, this makes it difficult to keep up with all the developments that are happening within the software and the scientific community. To address this growth in functionality, here we have compiled the latest method developments in the Rosetta software suite from the past five years, divided into several modeling categories as outlined below. This report is intended to serve as a guide for Rosetta users and developers — whether new, returning, or seasoned — who want to be updated on newest developments. The supplement contains a tour of the protocols with extensive links to documentation and resources on the web.

## Figure 1: Capabilities of the Rosetta macromolecular modeling suite

Popular tasks that can be addressed in Rosetta (blue) and major systems that can be modeled (red).



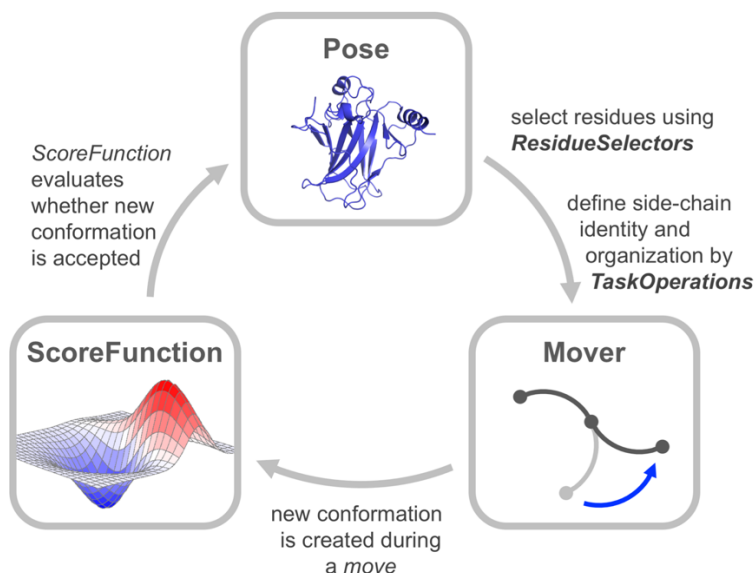
## 1. Major applications

The general outline of a typical Rosetta protocol is depicted in Figure 2: the conformation of a biomolecule (the *Pose*) is altered, either deterministically or stochastically, via a *Mover* and the resulting conformation is evaluated by a *ScoreFunction*. The *move* is accepted based on the Metropolis criterion

and the energy difference between the original and the new conformation. Many trajectories are generated, and the final models are evaluated based on the scientific objective.

### Figure 2: Main elements of a Rosetta protocol

Three main elements are required in a Rosetta protocol. The *Pose* is the biomolecule, such as a protein, RNA, DNA, small molecule, or glycan, in a specific conformation. Residues in the *Pose* can be selected via *ResidueSelectors* and the behavior for side-chain optimization or mutation can be defined by *TaskOperations*. Specific *Movers* then control how the conformation of the *Pose* is changed, and the new conformation is subsequently evaluated by a *ScoreFunction*. The Metropolis criterion decides whether the new conformation is accepted in the sampling trajectory. Many independent sampling trajectories are generated, and the final models are evaluated based on the purpose of the protocol.



### Predicting protein structures

Rosetta was originally developed for *de novo* protein structure prediction, which is accomplished by assembling fragments from known protein structures via a Monte Carlo procedure and scoring the models with an advanced scorefunction that balances physical and statistical potentials. Since optimizing the fragments for structure prediction can improve model quality, the original fragment picker application was re-implemented as an object-oriented framework that is vastly more flexible and allows incorporation of various types of restraints from secondary structure prediction or experimental data, for instance from NMR chemical shifts<sup>9</sup>. Improvements in homology modeling were achieved by multi-template modeling in RosettaCM, which combines the most homologous portions from multiple templates into a single model (called hybridizing) while modeling missing residues *de novo*<sup>10</sup>. If a template is absent, protein structures can be predicted *de novo*, which remains one of the most challenging tasks in structural biology, even though the incorporation of evolutionary coupling constraints (for instance from GREMLIN<sup>11</sup>) has led to enormous improvements in model quality. To thoroughly search the conformational space, an iterative hybridize approach was implemented. It uses a genetic algorithm that recombines models from an input pool to create models that have features from their parents but are also distinctly different. Creating several child models in each iteration, updating the input pool, and performing 30-50 iterations lead to improved model accuracy because features that are scored favorably by the scorefunction are repeatedly used in the recombination, such that the models in the pool converge over time. This approach has been used to improve model quality of *de novo* predicted models<sup>12</sup> as well as homology models<sup>13</sup>. Model refinement or generating ensembles of structures (useful in particular for design) can be accomplished by several algorithms in Rosetta: *FastRelax*<sup>14</sup>, *Backrub*<sup>15</sup>, or using vicinity sampling in the *KIC/Next-Generation-KIC* loop modeling algorithms<sup>16,17</sup> (see loop modeling section).



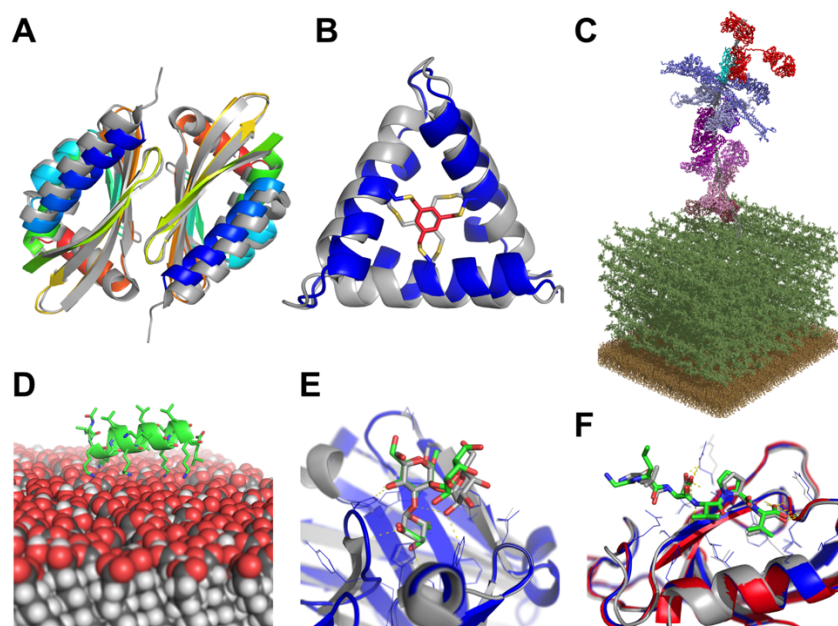
Experimental protein structure determination is challenging for proteins on solid surfaces such as biominerals, self-assembled monolayers, inorganic catalysts, and nanomaterials. RosettaSurface<sup>18</sup> samples protein conformations *ab initio* in both the solution and adsorbed states (Figure 3D) in order to account for adsorption-induced conformational changes. Experimental data can be incorporated into the simulation<sup>19</sup> to improve scoring, which remains difficult because the Rosetta scorefunction has been optimized for soluble proteins in aqueous solvent.

### Modeling protein–protein complexes

Another early Rosetta method was RosettaDock, which predicts the structure of protein-protein complexes from input monomers. The most recent iteration, RosettaDock4.0<sup>20</sup> incorporates protein flexibility from pre-generated protein ensembles, mimicking conformer selection. The new protocol has improved sampling efficiency by automatically adjusting the sampling rate based on the diversity of the input ensembles. Scoring has also been improved by using a novel, six-dimensional coarse-grained scoring scheme called *motif\_dock\_score*, which employs score grids generated from known complexes in the Protein Data Bank (PDB). In local docking benchmarks and backbone deviations of up to 2.2 Å, RosettaDock4.0 was able to successfully dock ~50% of complexes. For symmetric homomers, Rosetta SymDock2<sup>21</sup> can be used, which uses the same six-dimensional scoring scheme as in RosettaDock. Symmetry information can be extracted from a homologous complex, or a global docking search can be performed for a given point symmetry using Rosetta's symmetry framework<sup>22</sup>. An induced-fit based all-atom refinement relieves clashes in tightly-packed complexes to give physically realistic models. On a benchmark set of 43 complexes with different cyclic and dihedral symmetries, global docking on homology models had accuracies of 61% and 42% for cyclic and dihedral symmetries, respectively. These accuracies are substantially higher than for other symmetric docking tools and can be dramatically improved when adding restraints.

### Figure 3: Rosetta can successfully address diverse biological questions

(A) Overlay of the designed homo-dimeric curved  $\beta$ -sheet (dcs-E\_4\_dim\_cav3) in rainbow and the crystal structure in gray (PDBID 5u35). The protein is completely designed *de novo* and features a curved  $\beta$ -sheet, a large pocket, and a homodimer interface<sup>23</sup>. (B) Overlay of the *de novo* designed macrocycle 3H1 in blue and the NMR structure in gray (PDBID 5v2g). This is an example of a "CovCore" (covalent core) miniprotein that is held together covalently by a hydrophobic cross-linker at its core (in red for the design and gray for the NMR structure)<sup>24</sup>. (C) The interactome of M1 protein (the most important virulence factor of Group A *streptococcus*) and 15 human plasma proteins on the surface of bacteria (peptidoglycan layer (dark green), and the membrane (brown)). This 1.8MDa structure is measured directly in a complex mixture of intact bacteria and human plasma by PyTXMS. All structural models are provided by Rosetta where it consists of M1 protein (gray), IgG (red), four fibrinogens (dark to light blue), six albumins (dark to light pink), coagulation factor XIII A [F13A] (purple), C4bPa (cyan), haptoglobin [HP] (brown), and alpha-1-antitrypsin [Serpina1] (plum). This complex is supported by more than 200 chemical cross-links<sup>25</sup>. (D) Model of an LK- $\alpha$  peptide (LKKLLKLLKLLKL with a periodicity of 3.5 assuming a helical conformation) on a hydrophilic self-assembled monolayer surface. In solution, the peptide is unstructured, whereas when on the surface, experiments show that it assumes helical structure<sup>19</sup>. (E) Flexible docking of a carbohydrate antigen to an antibody. The crystal structure is in gray (PDBID 1mfa) and the Rosetta model in blue, with the carbohydrate in green. The coordinates for the antibody were taken from the PDB and the glycan coordinates started from a randomize backbone conformation and rigid-body orientation<sup>26</sup>. (F) High-resolution model of a peptide-protein complex generated using PIPER-FlexPepDock (model: blue; solved structure in gray, PDBID 1mfg). The predicted model was generated starting from a peptide sequence (LDVPV, derived from the C-terminal tail of ErbB2R) and the unbound structure of the receptor (Erbin PDZ domain, PDBID 2h3l, colored in red)<sup>27</sup>.



### Elucidating interfaces between proteins and small molecule ligands

Structure-based drug design has become a common approach for drug optimization due to increasing numbers of deposited structures in the PDB. RosettaLigand<sup>28</sup> has demonstrated success in predicting small molecule-protein interactions. Recent improvements to the algorithm rely on a low-resolution sampling step via the *TransformMover*, which combines translational and rotational perturbations in a single step, and using scoring grids for energy evaluation<sup>29</sup>. Further, the algorithm allows backbone flexibility, mimicking the induced fit hypothesis<sup>30</sup>. On a benchmark of 43 complexes, this new algorithm demonstrated an enhanced docking success by 10-15% with an effective 30-fold speedup over the original RosettaLigand performance, enabling virtual high-throughput screening (vHTS) of medium-sized ligand libraries. Later in the drug development process, when medicinal chemists optimize ligands based on structure-activity relationships (SAR), they synthesize ligands that typically share a core chemical scaffold and are assumed to bind to their target in a similar fashion<sup>31</sup>. RosettaLigandEnsemble<sup>32</sup> improves sampling during ligand docking by taking advantage of ligand similarities and docking a congeneric series of ligands at the same time, allowing for a placement that works for all considered ligands while optimizing the binding interface for each ligand independently. Experimental SARs can be included by promoting certain binding modes.

Another approach for therapeutic intervention is to use small molecule ligands as competitive inhibitors of protein-protein interactions. A common challenge, however, is that the protein's inhibitor-bound conformation often differs from the unbound or protein-protein bound conformation. The pocket optimization approach in Rosetta identifies protein surface pockets and uses their volume as an additional scoring term: this allows the user to start from an unbound protein structure and carry out biased sampling of a protein such that low-energy pocket-containing states are preferentially explored<sup>33,34</sup>. The specific conformations sampled through this approach match "druggable" alternate conformations observed in ligand-bound structures<sup>33,34</sup>, implying that these states are excellent starting points for virtual screening. The pockets sampled on the protein surface can then be matched to complementary ligands directly, by using the pocket itself as the starting point for pharmacophore-based screening<sup>35</sup>. Alternatively, these pockets can also be used for Rosetta's *Docking Approach using Ray Casting* (DARC<sup>36</sup>) method. DARC uses ray-casting to rapidly position a ligand in the protein surface pocket<sup>36</sup>; by iterating over many candidates, DARC provides a means for very rapid virtual screening. DARC has also been adapted for GPUs<sup>37</sup>, and the newer implementation<sup>38</sup> includes features that provide improved performance in virtual screening benchmarks.

### Designing new proteins and functions

The inverse problem of protein structure prediction is protein design, where the objective is identification of a sequence that best represents a given structure. In particular, *de novo* design and design of novel protein functions towards therapeutic intervention remains one of the grand challenges in structural biology. This problem is addressed by various methods in Rosetta. The SEWING protocol creates *de novo* designs by recombining parts of protein structures from randomly-selected helical building blocks<sup>39</sup>. SEWING's requirement-driven approach allows users to specify features or properties that should be incorporated into their designs during backbone generation without necessarily requiring a certain size or three-dimensional fold. New features include incorporation of functional motifs such as protein-binding peptides for protein interface design and partial or complete ligand binding sites for ligand-binding protein design<sup>40</sup>. A somewhat similar algorithm has been implemented for antibody design (AbDesign, see "modeling antibodies" below), which was generalized for enzyme design<sup>41</sup>. A more general approach is RosettaRemodel, which performs protein design by rebuilding parts or all of the structure<sup>42</sup> from fragments of known proteins structures. It relies on a blueprint file in which the user defines secondary and supersecondary structure of the fold to be built. Remodel interfaces with a number of Rosetta protocols and can be used for various applications such as *de novo* modeling, fixed-backbone sequence design, refinement, loop insertion, deletion, and remodeling, as well as disulfide engineering, domain assembly, and motif grafting.

For designing multifunctional proteins such as biosensors, bioswitches and tunable affinity clamps, Blacklock et al. developed the domain insertion protocol LooDo (Loop-directed domain insertion). With LooDo, proteins are designed by inserting a domain into another by two flanking linker regions<sup>43</sup>. The linker regions are sampled *via* fragment insertion to determine relative positioning of the domains, followed by generalized kinematic loop closure<sup>44</sup> (GenKIC, see below) and enzyme design to optimize the interface.

In protein design, a common task is not only design towards a certain goal (positive design), but additionally, design away from undesired features (negative design). Such a Multi-State Design<sup>45</sup> (MSD) approach evaluates strengths and weaknesses of a single sequence on multiple backbones to account for positive and negative design, for instance binding to one but not another protein partner. REstrained CONvergence<sup>46</sup> (RECON) takes this idea one step further by allowing each state to sample multiple sequences during the design process, which is iteratively applied by increasing the restraint weight to encourage sequence convergence. RECON is effective for large multi-state design problems, such as antibody affinity maturation or prediction of evolutionary sequence profiles of flexible backbones<sup>47,48</sup>.

*De novo* protein design is somewhat easier for structures that are primarily highly regular helices and sheets, as the principles dictating their conformations are well known. However, designing curved and twisted  $\beta$ -sheets requires a deeper understanding of the structural irregularities that enable them. These principles were implemented in the curved  $\beta$ -sheet design method to design a variety of protein folds with curved sheets (Figure 3A), creating pockets suitable for tailoring ligand-binding and enzymatic active sites<sup>23</sup>.

During computational *de novo* protein design<sup>49</sup>, a stringent test for the consistency of the designed sequence is whether *ab initio* structure prediction will yield the same structure that was used as a starting point for the design. However, computationally testing a large number of designs is prohibited by the vast conformational search space for *ab initio* structure prediction. To drastically limit that search space and test many more designs, the biased forward folding method<sup>23</sup> uses three (instead of the typical 200) fragments per residue position. Fragments are chosen based on the RMSD to the native structure in design. The designs achieving near-native sampling are more likely to have funnel-shaped energy landscapes and therefore worth assessing with *ab initio* structure prediction.

Design of protein function has been accomplished by grafting a known motif from a template protein onto a new protein (*motif grafting*). This approach has been used for antibodies and for vaccine design<sup>50</sup> using the *fold\_from\_loops* application, where the functional motif is used as a starting point of an extended structure that is folded following the constraints of a target topology. Iterative refinement is carried out via sequence design and structural relaxation before filtering and human-guided optimization. This



application has been extended into the *Functional Folding and Design* (FunFolDes) protocol, which includes multi-segment motif grafting, different residue length motif insertion, incorporation of restraints, and folding in the presence of a binding target<sup>51</sup>. Fragments selected according to the structure of the target topology improve the performance of the folding stage via the *StructFragmentMover*.

### Designing interfaces between proteins and interaction partners

Problems related to protein design include designing interfaces of proteins with their interaction partners such as proteins or small molecule ligands, and predicting  $\Delta\Delta G$ s of mutation (e.g. alanine scanning). Predicting  $\Delta\Delta G$ s of mutations for protein stability or protein-protein interactions is a difficult problem with low correlation coefficients (0.5-0.7)<sup>52</sup>, because the effect of the mutation is small compared to the total energy in the system, and because protein flexibility adds noise to the energies that can mask the effect of mutations. In the simplest case of alanine scanning (mutating into Ala), methods that use a “soft-repulsive” energy function without modeling backbone flexibility<sup>53,54</sup> have typically outperformed methods that allow protein flexibility and use hard-repulsive energy functions<sup>55</sup>. FlexDDG<sup>56</sup> was created to improve protein-protein interface  $\Delta\Delta G$  predictions and generalize them to residues other than Ala. The protocol creates conformational ensembles using Rosetta backrub sampling<sup>57</sup>, then repacks sidechains, minimizes torsions and computes change in protein-protein interaction  $\Delta\Delta G$  by averaging across the ensembles. On 1240 interface mutants, FlexDDG outperforms Rosetta’s *ddg\_monomer* application, which was originally created and validated to predict changes in stability upon mutation, not interfaces.

Designing ligand-binding interfaces in proteins is challenging due to inaccuracies in the energy function (and implicit solvation), the flexibility of ligands, and the sensitivity of protein-ligand interactions to even subtle conformational changes<sup>58</sup>. Flexible backbone design methods that use pre-generated ensembles as a starting point for design<sup>59,60</sup> perform poorly in benchmarks, likely because the ensemble does not accurately describe the unbound-to-bound conformational changes. The *CoupledMoves* protocol couples backbone flexibility with changes in sidechain rotamers or ligand orientation or conformers, and leads to substantial improvements in various benchmarks<sup>61</sup>.

Symmetric protein assemblies can now be modeled using parametric design. Nature created super-helical coiled-coils that are well-described by geometric equations using Crick parameters<sup>62</sup>, which include variables for the radius of the bundle, major helical twist, minor helix rotation about the primary axis, etc. Several Movers such as *MakeBundle*, *PerturbBundle*, and *BundleGridSampler* allow designing helical bundles<sup>24,63</sup> and  $\beta$ -barrels based on pre-defined or sampled parameters. Since parametric methods do not rely on fragments libraries, these modules can be applied to non-canonical coiled-coil heteropolymers.

### Modeling peptides and peptidomimetics

The inherent flexibility of peptides imparts a large conformational search space to them, which leads to challenging modeling problems; when peptide modeling is combined with another simulation, e.g. docking, the increase in conformational space makes the modeling task virtually impossible using traditional approaches. FlexPepDock addresses this problem by allowing targeted sampling of the peptide flexibility during its docking into a given binding site, either by refining an approximate peptide conformation (FlexPepDock refinement<sup>64</sup>), or by full *ab initio* sampling of the peptide conformation (FlexPepDock *ab initio*<sup>65</sup>). Peptide docking is especially challenging when the binding site on the receptor is unknown. However, it can be simplified based on the observation that (for peptides built from canonical amino acids) the bound peptide conformation is often included in the fragments generated by the *FragmentPicker*. The PIPER-FlexPepDock<sup>27</sup> protocol rigid-body docks these fragments using PIPER FFT-based docking<sup>66</sup>, and refines the complex using FlexPepDock<sup>64</sup>. PIPER-FlexPepDock can generate highly accurate peptide-protein complexes from a peptide sequence and a free receptor structure (Figure 3F).

Many protein-protein interactions (PPI) are mediated by often disordered peptide segments that are responsible for most of the binding energy<sup>67–70</sup>. PeptiDerive<sup>71</sup> detects such segments in a PPI complex through a sliding window approach. PeptiDerive was extended to cyclized peptides and is available on the ROSIE<sup>72</sup> server.

Conformations of cyclic peptides can be sampled with *simple\_cycpep\_predict*, which restricts the conformational search space through cyclization<sup>24,44,73</sup> via the Generalized Kinematic Closure (GenKIC) algorithm (see “loop modeling” below). *Simple\_cycpep\_predict* does not rely on protein fragments and can model non-canonical chemistries (Figure 3B), being a generalization of earlier protocols.

Multi-specificity is common at protein-peptide interfaces, meaning that the protein can interact with multiple substrates at the same interaction site. This can be exploited for identifying and designing novel substrates. Multi-specificity can be modeled with MFPred<sup>74</sup>, which is a rapid, flexible-backbone self-consistent mean field theory-based technique. MFPred can predict experimentally determined peptide specificity profiles for a range of receptors, at equivalent or better prediction accuracy and a 10- to 1000-fold lower computational cost when compared to other methods.

**Loop modeling for structure prediction and design**

Loop modeling was implemented early in Rosetta<sup>75,76</sup> to generate structures for loops or gaps in models, with initial approaches relying on fragments to sample conformations and the iterative Cyclic Coordinate Descent (CCD) algorithm<sup>77</sup> for chain closure. Subsequent developments introduced inverse kinematic closure methods (termed “KIC”) into Rosetta, that rely on polynomial resultants to analytically solve for closed conformations and which produced more native-like loop conformations<sup>78</sup>. KIC was used for modeling protein surface and interface loops, and to design and refine active site loops or regions binding small molecules<sup>79</sup>. Next-Generation KIC (NGK)<sup>17</sup> made improvements to sampling loop conformations by employing diversification (i.e. sampling a wider range of possible conformations) and intensification (i.e. to focus sampling on previously generated conformations) to identify lowest-energy conformations. NGK substantially increases the fraction of near-native models<sup>17</sup> and allows modeling longer loops. GeneralizedKIC<sup>44</sup> (GenKIC) samples or perturbs loop geometries between fixed endpoints for any continuous peptide chain including those with non-standard peptide chemistries, for instance for non-canonical backbones. GenKIC can sample backbone conformations containing L- and D- $\alpha$ -amino acids,  $\beta$ -amino acids, peptoids, oligoureas, or more exotic chemical building-blocks for which template structures do not exist in structural databases. Additionally, GenKIC can sample conformations involving side-chain connections (e.g. disulfide bonds, side-chain isopeptide bonds, etc.), covalently-attached ligands or crosslinkers, or chemistries that conventional loop-modelling algorithms do not typically handle.

Most Rosetta loop modeling algorithms were primarily developed for structure prediction. However, design constitutes the opposite problem, finding low-energy sequence–structure combinations that satisfy certain design goals. LoopHashKIC<sup>80</sup> addresses this problem and uses the Rosetta LoopHash algorithm<sup>81</sup> to efficiently query a database of loop conformations based on rigid-body transforms between the first and last loop residues. LoopHashKIC uses LoopHash to identify a suitable peptide fragment, and then uses KIC to find an exact solution to close the backbone. To improve the local sequence-structure compatibility in *de novo* designed loops, the *ConsensusLoopDesign* task operation accessed through Rosetta Scripts allows a user to restrict the amino acid identities of loops based on sequence profiles of naturally occurring loops with the same region of backbone dihedral angle space (Ramachandran bins)<sup>23,82</sup>.

**Modeling antibodies and other proteins in the immune system**

Due to the therapeutic significance of antibodies, several protocols have been developed in our community for structure prediction, docking and design that involve antibodies and other proteins in the immune system, such as T-cell receptors (TCR), displayed antigens of the Major Histocompatibility Complex (MHC) and other soluble antigens and immunogens. An overview of the various applications can be found in Table 1.

protocol	target	task	comments
RosettaAntibody <sup>83–86</sup> Gray lab	antibodies	homology modeling / docking	Identifies templates, assembles them into a structure, and models loops <i>de novo</i> while refining VH-VL orientation <sup>87</sup> ; allows multiple templates; uses a key constraint <sup>88,89</sup> for CDR H3 modeling; good for modeling camelid antibodies <sup>90</sup> and antibodies on the scale of the human repertoire <sup>91,92</sup>
AbPredict <sup>93</sup> Fleishman lab	antibodies	structure prediction	Does not rely on templates, samples backbone fragments and rigid-body orientations from known structures, without

			considering sequence homology, hence being able to model antibodies accurately with sequence identity as low as 10%; AbPredict2 available as webserver <sup>94</sup>
RosettaMHC <sup>95</sup> Sgourakis lab	antigen / MHC-I / (chaperone or T-cell receptor)	modeling / docking	Predicts peptide antigens that bind to all known MHC-I alleles models peptide/MHC-I structures <sup>96</sup> ; code implementation in PyRosetta
RosettaTCR <sup>97</sup> Pierce lab	T-cell receptors	structure prediction	Models TCRs from sequence, via template identification, grafting of loop templates onto framework regions, minimization and loop refinement; to gain structural insights into TCRs, e.g. those targeting cancer neoepitopes <sup>98</sup> , or to identify features of sets of TCRs from high throughput sequencing; can be combined with docking to generate models of TCR-peptide-MHC complexes <sup>99</sup> or TCRs in complex with non-peptide antigens bound to MHC-like proteins <sup>100</sup> .
SnugDock <sup>101</sup> Gray lab	antibody-antigen	docking	Input is starting conformation and optionally an ensemble of antibodies/antigens, then runs local docking to refine both the antibody-antigen interface and the heavy-light chain interface (within the antibody) and re-models the CDR H2/H3 loops at the interface; also includes a CDR H3 structural constraint <sup>88,89</sup> and docking of camelid antibodies <sup>90</sup>
RosettaAntibodyDesign <sup>102</sup> (RABD) Dunbrack lab	antibody-antigen	design	Based on RosettaAntibody <sup>85</sup> ; design of specific CDRs of different clusters and lengths, sequence design using cluster-based CDR profiles or conservative mutations, or <i>de novo</i> design of whole antibodies; based on the North-Dunbrack CDR clustering <sup>103</sup> , reducing deleterious sequence mutations; benchmarked on a diverse set of 60 interfaces from both lambda and kappa antibodies; experimental benchmarking of two complexes showed affinity improvements between 10 and 50-fold
Epitope removal <sup>104,105</sup> Baker lab	MHC epitopes	design	Includes experimental immunogenic epitope data, MHC epitope prediction tools, and host genomic data to enable protein design with reduced immunogenicity while retaining function and stability <sup>104</sup> ; uses a machine learning-based epitope prediction for 28 different alleles, restricts design to select 15mer epitope regions, and uses a greedy stepwise protein design <sup>105</sup> to eliminate the most immunogenic epitopes with as few mutations as possible, avoiding disruptive core mutations likely to destabilize the protein
AbDesign <sup>106,107</sup> Fleishman lab	antibodies	design	Cuts experimentally determined antibody structures along conserved positions to create interchangeable segments, recombines them to produce novel antibody models <sup>106,107</sup> ; models are docked to a target of interest, either locally to a specific epitope, or globally, followed by optimization comprised of backbone sampling and sequence design for improving stability and binding affinity

### Using experimental data to direct modeling

The use of experimental data in modeling can vastly restrict the conformational search space, therefore allowing the modeling of larger, more complex biomolecules to greater accuracy. Electron density maps from cryo-electron microscopy (cryoEM) or X-ray crystallography have become more readily available in the past decade and methods to incorporate these types of data have been successfully used for high-resolution structure determination. Since cryoEM density maps are often of low resolution, *de novo* structure determination methods require a combinatorial search procedure to unambiguously assign all densities to residues in the protein. A *de novo* method described by Wang et al. applies a model building approach<sup>108</sup> for density maps between 3-5Å that fits fragments into densities and scores their match based on secondary structure, fit with density, loop closure, clashes and consistency between

overlapping fragments to assign sequence into densities. While this method requires >70% of the map to be assigned initially, an updated version of this method, the RosettaES<sup>109</sup> enumerative sampling approach, forgoes this requirement. RosettaES gradually extends the model one residue at a time until all residues have been assigned. At each iteration, short fragments are used to sample the nearby conformational space of the growing model, while undergoing a series of clustering and filtering steps based on the Rosetta energy and fit to the density.

If assignment is complete but the data are low-resolution, refinement into density maps is necessary. Several methods have been developed for density maps in the 3.0-4.5Å resolution range. One method<sup>110</sup> iterates between refinement with Phenix in reciprocal space to physically plausible conformations, and Rosetta in real space, because Rosetta's all-atom scorefunction compensates for the lack of high-resolution data, while the density map restrains backbone and side-chain sampling in real space. Refinement can also be seeded from homology models, followed by density-guided rebuilding and refinement of coordinates and B-factors<sup>111</sup>. More recently, an automated fragment-guided refinement pipeline<sup>112</sup> splits the density map into independent training and validation maps. It finds regions with poor density fit, iteratively rebuilds them with fragments using the training map, filters the models based on their fit to the validation map, model geometry from MolProbity and fit to the full map, and then optimizes against the full map. The frameworks for electron density maps and carbohydrate modeling<sup>26</sup> in Rosetta (below) were connected<sup>113</sup> for refinement of carbohydrates into low-resolution electron density maps from cryoEM or crystallography.

NMR data were incorporated into *de novo* structure prediction early in Rosetta development, creating RosettaNMR. Chemical shifts were used for fragment picking using CS-Rosetta<sup>114</sup>, which could be used in conjunction with NOE, RDC<sup>115</sup>, PCS<sup>116–118</sup> and PRE data. Improvements, for instance through RASREC resampling<sup>119</sup> allowed the use of sparse<sup>120</sup> or unassigned data<sup>121</sup>, easier-to-obtain data (backbone-only<sup>122</sup>), modeling larger and more complex proteins<sup>123</sup>, membrane proteins<sup>124</sup>, symmetric systems<sup>125</sup>, and combination with data from SAXS<sup>126</sup>, cryoEM<sup>127</sup>, distance restraints from homologous proteins<sup>128</sup> and evolutionary couplings<sup>129</sup>. CS-Rosetta also has the AutoNOE<sup>130,131</sup> module for automatic assignment of NOESY data for use in structure calculations. RosettaNMR was recently overhauled and reconciled with CS-Rosetta and PCS-Rosetta to allow seamless integration of several types of NMR restraints (CS, RDC, PCS, PRE, NOE) in one consistent framework<sup>132</sup> that could be applied to structure prediction, protein-protein docking, protein-ligand docking, and symmetric assemblies.

Covalent labeling mass spectrometry data provides information on relative solvent exposure of residues, therefore yielding information on protein tertiary structure. A low-resolution score term from hydroxyl radical foot-printing has been implemented in Rosetta that can improve model quality in structure prediction<sup>133,134</sup>. Finally, data from chemical cross-linking mass spectrometry has been incorporated into an automated workflow to identify protein-protein interactions. The PyTXMS<sup>25</sup> method combines the sensitivity of mass spectrometry to analyze complex samples with the power of Rosetta structural modeling and protein-protein docking to efficiently sample the vast conformational space and identify interactions (Figure 3C). A machine learning algorithm based on high resolution MS1 data guide the potential binding interface selection which is then validated and adjusted by a repository of structural models and MS2 (DDA) samples.

### Modeling nucleic acids and their interactions with proteins

Several advances have been made in the representation of nucleic acids in Rosetta. The *StepWise Monte Carlo* protocol (SWM) has achieved RNA structure predictions reaching atomic accuracy<sup>135</sup>; the approach provides an acceleration over the original enumerative *StepWise Assembly* (SWA) method<sup>136,137</sup>. A version of SWA that rebuilds one nucleotide at a time enables fine-grained correction of errors in RNA coordinates fit into crystallographic or cryo-EM maps by *Enumerative Real-space Refinement ASsisted by Electron density under Rosetta*<sup>138,139</sup> (ERRASER).

The *Fragment Assembly of RNA with Full-Atom Refinement* (FARFAR) structure prediction protocol<sup>140,141</sup> also permits working with chemically modified nucleotides, picking fragments for the most chemically similar base available<sup>135</sup>. Homologous fragments can automatically be eliminated from fragment sets to

give pseudo-blind prediction results<sup>135</sup>. As another connection of Rosetta's RNA tools to experimental structural biology, <sup>1</sup>H NMR data can be used for RNA modeling via the CS-Rosetta-RNA protocol<sup>142</sup>.

The most recent advances in Rosetta RNA tools expand the fragment assembly protocol to support modeling RNA-protein complexes through simultaneous folding and docking<sup>143</sup>. RNA-protein interactions are handled via additional knowledge-based score terms that supplement the low-resolution RNA scorefunction. Free energy perturbations from RNA or protein mutations can be modeled with the Rosetta-Vienna  $\Delta\Delta G$  protocol<sup>144</sup>. Structure coordinates can further be built into cryo-EM density maps for large RNA-protein complexes with DRRAFTER (*De novo Ribonucleoprotein modeling in Real space through Assembly of Fragments Together with Experimental density in Rosetta*)<sup>145</sup>.

Redesign and prediction of protein-DNA interfaces is also possible in Rosetta<sup>146,147</sup> and has been accomplished with flexible protein backbones<sup>148</sup>, genetic algorithms<sup>146,148,149</sup> and motif-biased rotamer sampling<sup>150,151</sup>. However, the biggest limitation of these approaches is that they rely on fixed DNA backbone conformations, which in nature can be highly flexible. Key to successful protein-DNA design is an energy function that is optimized<sup>151,152</sup> for these highly polar and solvated interfaces. Rosetta further supports prediction of specificity and affinity<sup>153</sup> and the prediction of DNA binding preferences of homologous proteins. Multi-template modeling in RosettaCM<sup>154</sup> was successfully applied to this challenge<sup>155</sup>. To accomplish this, protein homology modeling was followed by docking of multiple competing DNA sequences threaded onto the original crystal structure backbone and comparing the energies of the resulting protein-DNA complexes.

### Modeling membrane proteins

Membrane proteins constitute about 30% of all proteins and are targets for over 60% of pharmaceuticals on the market<sup>156</sup>. However, experimental difficulties have limited our understanding of their structures<sup>157</sup>. Previously, Yarov-Yarovoy<sup>158,159</sup> and Barth<sup>160</sup> implemented tools for low- and high-resolution structure prediction of membrane proteins, termed RosettaMembrane. These tools were recently re-engineered for compatibility with Rosetta3<sup>2</sup> into a platform called RosettaMP<sup>161</sup>. RosettaMP implements core modules for representing, sampling, and scoring proteins in the context of an implicit membrane. RosettaMP is compatible with key modeling protocols including docking, design,  $\Delta\Delta G$  prediction<sup>52</sup>, PyMOL visualization<sup>162</sup>, and assembly of symmetric proteins. In addition, a set of basic modeling tools<sup>163</sup> is implemented, for instance for scoring, transforming a membrane protein into the membrane coordinate frame, *de novo* modeling for single transmembrane span helices, introducing mutations, and visualization in the membrane. RosettaMP has further enabled rapid development of new modeling tools including structure-based detection of lipid exposed residues in the membrane<sup>164</sup> environment and domain assembly of full-length protein models from structures of transmembrane and soluble domains<sup>165</sup>. The RosettaCM protocol for multi-template homology modeling has also been adopted for membrane proteins<sup>8</sup>.

### Adding carbohydrates to the modeling process

Carbohydrates are fundamental to life<sup>166,167</sup>, but because of challenges in experimental characterization and computational sampling and scoring, their structures have been historically under-studied. The RosettaCarbohydrate framework<sup>26</sup> allows modeling of carbohydrate structures and complexes. The framework is integrated into Rosetta such that it is possible to model glycosylated proteins or protein-sugar complexes (Figure 3F) with the same algorithms one would use for proteins. RosettaCarbohydrate is not limited to commonly studied sugars but can handle the full gamut of carbohydrate structures, including linear, cyclic, and branched structures, sugar modifications, and conjugations. Methods exist for sampling ring conformations, packing substituents, refining glycosidic linkages, sampling from linkage "fragments", and extending glycan chains. Scoring of saccharide-containing sugars includes a quantum-mechanically derived intrinsic backbone term<sup>168</sup>. Because saccharide residues are stored as distinct data structures, we can integrate bioinformatic and statistical data into our algorithms, which opens the doors for glycoengineering and design applications. RosettaCarbohydrate has been integrated with various other frameworks in Rosetta, such as loop modeling (GenKIC and Stepwise Assembly), refinement (*GlycanTreeModeler*), symmetry, and RosettaScripts-accessible classes such as *MoveMaps* and



*ResidueSelectors*. Linkages are automatically determined during PDB read-in. Carbohydrates work with Cartesian minimization, and they can be refined into electron density maps<sup>113</sup>.

## 2. The brain of Rosetta: its scorefunction

Rosetta's energy function has been continuously improved over many years<sup>169</sup> and a full review of the modern all-atom energy function was published recently<sup>3</sup>. Briefly, the newest energy function REF2015<sup>4</sup> (REF stands for Rosetta Energy Function) features two main advancements. First, reproducibility of thermodynamic observables (such as liquid-phase properties<sup>170</sup> and liquid-to-vapor transfer free energies<sup>171</sup>) was added to the optimization objectives, in addition to structure<sup>172</sup>-based tests. Second, a new, derivative-free optimization technique was developed, which is suitable for robust optimization of >100 parameters. Further, a new energy term was added that takes into consideration non-ideality of bond lengths and angles in cartesian space<sup>173</sup>. The cartesian term<sup>173</sup> is also the basis for a *cartesian\_ddG* method that has been used to calculate  $\Delta\Delta G$ s of mutation to probe changes in protein stability. Only the backbones and side chains of residues nearby the mutation site are allowed to move<sup>174</sup>. Due to the local optimization, this protocol is much faster than *ddg\_monomer*<sup>175</sup>, while retaining the same level of accuracy. The default Rosetta energy function is now also compatible with an expanded palette of chemical building-blocks: canonical and non-canonical L- $\alpha$ -amino acids and their D-amino acid counterparts, exotic achiral amino acids like 2-aminoisobutyric acid (AIB), peptoids, and oligoureas. The ability to model metalloproteins has also been added to Rosetta<sup>176,177</sup>. As noted above, scorefunctions that enable simultaneous modeling of protein and RNA are being explored<sup>144</sup>. The scorefunction is now thread-safe and fully mirror symmetric, i.e. enantiomers in mirror conformations score identically. Guidance energy terms for design have been added to encourage certain features, such as specific amino acid compositions<sup>44,73</sup>, hydrogen bonding networks, or global or local net charges, and discourage others, such as repeat sequences that hinder NMR assignments, buried unsatisfied hydrogen bond donors and acceptors, or voids within the protein<sup>178</sup>. Further, a consensus scoring method, which utilizes the semi-orthogonal nature of the Rosetta and Amber energy functions, was developed for model ranking to identify most near-native models<sup>179</sup> from the pool of generated decoys. This approach led to the development of a Python-based tool (AMBRose) for interconversion between Rosetta and Amber models to facilitate consensus scoring.

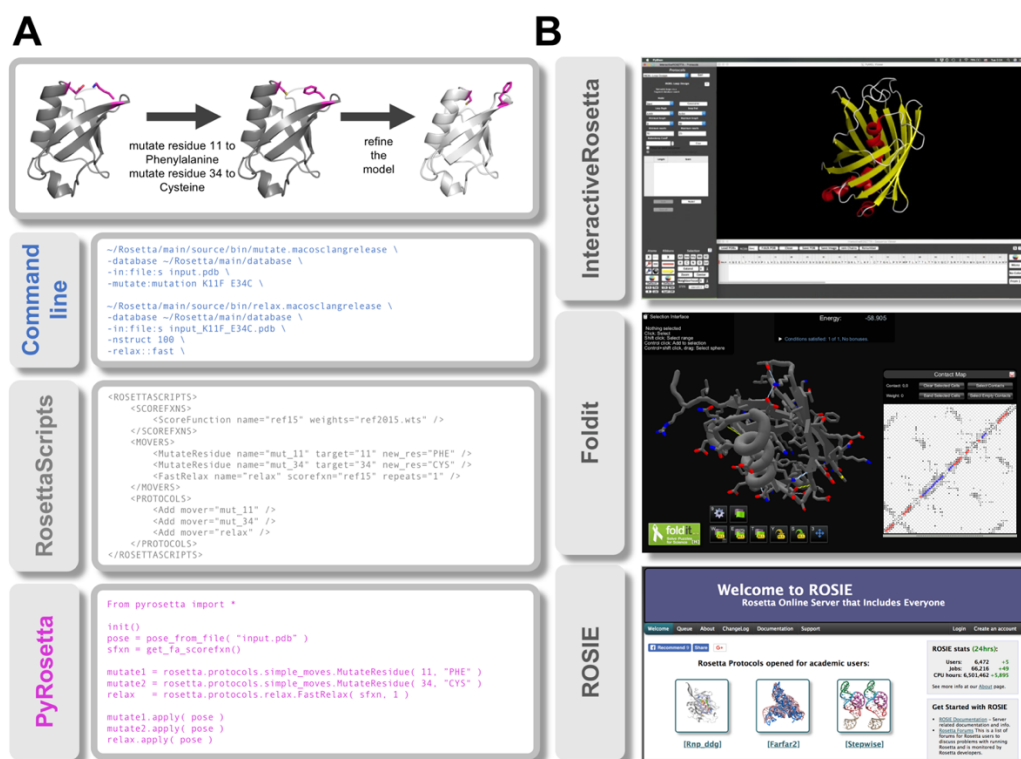
Hydrogen bond networks are important for biomolecular structure and catalysis but have been challenging to design because of pairwise interactions that have multi-body, cooperative properties. The HBNet protocol<sup>180</sup> has been used to design *de novo* coiled coils with interaction specificity mediated by designed hydrogen bond networks, including homo-oligomers<sup>180</sup>, membrane proteins<sup>63</sup>, and large sets of orthogonal heterodimers<sup>181</sup>. An improvement to HBNet uses a Monte Carlo search procedure to sample hydrogen bond networks with drastically improved performance<sup>182</sup>. We further developed a statistical potential to place highly-coordinated water molecules on the surface of biomolecules. On a data set of 153 high-resolution protein-protein interfaces, the method predicts 17% of native interface waters with 20% precision within 0.5 Å of the crystallographic water positions<sup>183</sup>. The potential is accessible through the WaterBoxMover in RosettaScripts.

## 3. User interfaces and usability

Advances in Rosetta have also focused on improving usability of the software through developing several user interfaces to suit different use cases and styles (Figure 4). The command line interface was the first and is still the most-often used interface to Rosetta methods. Structure input and output was enhanced by the ability to read and write mmCIF files (via an external library) using the same mechanisms as PDB files, which permits representation of large complexes that are ill-suited for the PDB format (e.g. the ribosome). This comes with the ability to read the Protein Databank's Chemical Component Dictionary, the description of the chemical composition of residues in the officially released PDB structures. Multithreading support has been added to Rosetta, which required a major refactor of its core architecture for thread-safety, allowing shared-memory parallelism. Multithreading is currently available for specific protocols (*simple\_cycpep\_predict*) with planned expansion to other applications (including the *JobDistributor* jd3).

### Figure 4: User interfaces to Rosetta

(A) Rosetta can be run from a terminal and offers three different interfaces to the codebase. The top panel outlines the task to be accomplished: making two mutations in a protein and then refining the structure. The panels underneath show how this task can be accomplished in the different interfaces. The command line panel shows the executable, input files and options to run two specific Rosetta applications. RosettaScripts is an XML-based scripting language that offers more flexibility by combining *Movers* and *ScoreFunctions* into a custom *Protocol*. PyRosetta offers direct access to the underlying Rosetta code objects but requires knowledge of the Rosetta codebase. (B) Point-and-click interfaces to Rosetta. InteractiveRosetta is a graphical user-interface (GUI) to PyRosetta. It offers controls to the most popular protocols, file formats and options. Foldit is a videogame primarily used to crowd-source real-world scientific puzzles but can also be used on custom proteins of interest. It allows access to some popular applications via a game interface. ROSIE is a super-server hosting a multitude of servers each executing a particular Rosetta protocol. It currently includes servers for 21 Rosetta methods. [The InteractiveRosetta panel was reproduced with permission from *Bioinformatics*.]



In addition to the command line interface, Rosetta features two major scripting interfaces: RosettaScripts and PyRosetta. RosettaScripts<sup>6</sup> is a popular scripting interface that uses Extensible Markup Language (XML) to build fairly complex protocols using core Rosetta machinery<sup>2</sup>. Comprehensive knowledge of the codebase is unnecessary since most of the underlying modules<sup>8</sup> have been thoroughly documented – documentation is now also generated using XML schema, which validates the RosettaScripts XML files at runtime. RosettaScripts was further extended and generalized to enhance consistency: *ResidueSelectors* enable selection of residues based on specific properties such as chain, amino acid, secondary structure, index, solvent accessible surface area, and others, and can be used in conjunction with *MoveMapFactories*, which control a structure's flexibility during energy minimization. *ResidueSelectors* are also accepted by *TaskOperations* which control side-chain identity and optimization. A more general analysis tool, *SimpleMetrics*, allows custom analyses of models through RosettaScripts and writes the output into the scorefile. The *SimpleMetrics* system is more integrated and robust than previous tools, such as the *InterfaceAnalyzer* or the *FeaturesReporter*.

PyRosetta<sup>5,184</sup> is a collection of Python bindings to the source code, exposing ~7,400 classes and 88,000 functions. PyRosetta allows custom protocol development that is flexible and fast, but it requires familiarity with the underlying structure of the codebase. Not all of options available in RosettaScripts have corresponding API-level configuration, so in order to take full advantage of those protocols, PyRosetta can now configure objects using RosettaScripts XML. This brings the added advantage of harmonizing the documentation across multiple interfaces.

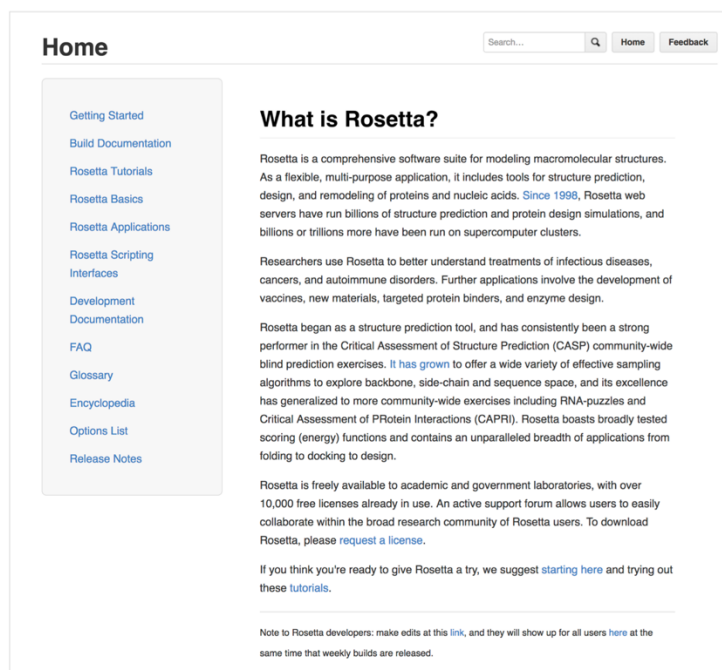
InteractiveROSETTA<sup>185</sup> is a graphical interface for PyRosetta that presents easy-to-use controls for several of the most widely-used Rosetta protocols alongside a selection system that uses PyMOL as a visualizer. InteractiveROSETTA is capable of interacting with remote servers running a standalone Rosetta install, rendering it easy to incorporate more sophisticated protocols that are not accessible in PyRosetta and/or require significant computational resources.

Foldit Standalone<sup>186,187</sup> is a graphical interface to Rosetta based on the Foldit video game<sup>7,188</sup>. Foldit Standalone provides several interactive structure manipulations, including pulling directly on the structure, rigid body docking, and residue mutation, insertion and deletion. Users can apply hard and soft constraints that guide automated moves such as packing and minimization, and provides real-time scoring updates as the structure changes. Additional features include multiple sequence alignments for template-based modeling, along with electron density-, Ramachandran-, and contact-map visualizations. Further, scientists and educators can now run their own custom Foldit puzzles for a group of their choosing, a new feature called “Custom Contests”<sup>189</sup>.

The Rosetta community has further devoted an enormous effort to enhance the user friendliness of Rosetta by rewriting and adding documentation (Figure 5). We now use a public-facing Gollum wiki (<https://www.rosettacommons.org/docs/latest/Home>) for various levels of documentation, such as application documentation, tutorials for beginning users, and static protocol captures that accompany manuscripts for scientific reproducibility (see supplement for links). The Gollum wiki is easily editable by members of the RosettaCommons. Command line executables accept a `-info` option, which prints relevant options for the current application in RosettaScripts, and debugging command lines is facilitated by improved error messages. Default, system-wide options (e.g. database paths) can now be specified in a *rosetta.rc* file. Lastly, code development in C++ is now easier with the help of available code templates that create much of the boilerplate code required to extend Rosetta.

### Figure 5: Main external documentation page for Rosetta

In 2015, our community performed a complete overhaul of the documentation for Rosetta. Documentation is now hosted on a Gollum wiki, which is version controlled and easily editable for members of our community. Accessibility and ability to edit the documentation has drastically improved the user-experience of Rosetta.



A limitation of Rosetta is the need for a local installation and compilation in a Unix-like environment. Webservers provide a user-friendly alternative and a number of independent servers have emerged in the Rosetta community. However, implementing and maintaining such servers comes at a substantial cost. To make it easier to provide Rosetta protocols as a webserver, ROSIE (Rosetta Online Server that Includes Everyone)<sup>72,190</sup> (<http://rosie.rosettacommons.org/>) provides a simple framework for “serverification” of protocols. ROSIE currently contains 21 webservers, with additional protocols continually being added.

### Conclusion

The Rosetta software is developed by a large, global community that aims to solve complex problems through collaborative solutions and code implementations. In the last five years, great strides have been made in Rosetta. More protocols are available now that enable modeling a broader range of biologically- and chemically-realistic systems, larger macromolecular complexes and in general more sophisticated systems. Prediction accuracies have improved through advances in the scorefunction, which is a combination of physics-based and knowledge-based potentials that were also fit against thermodynamic observables. Incorporating experimental data into the modeling process has also been facilitated and improved. Further, our community saw the need to develop more general, reusable, user-friendly, and scientifically reproducible protocols. This was motivated by the growth of the software and the developer community, the various user interfaces, the diversity of the community<sup>1</sup>, and the complexities of the protocols in this complex piece of software. The improvements to documentation allow users to quickly start using or developing custom protocols, while facilitating user support for the various Rosetta interfaces (command line, RosettaScripts, PyRosetta, etc.). Over the years, these applications have moved away from simply tackling basic science questions to more application-based scientific developments. The myriad of advances described above have made integration of Rosetta into existing experimental and computational scientific workflows increasingly useful and standard. Rosetta’s predictive powers can be used not only to analyze and verify existing data but to inform experiments to

galvanize engineering tomorrow's industrial enzymes, enable the creation of novel biomaterials, and accelerate the discovery of new potent, life-saving therapeutics that will change the world as we know it.

### **Acknowledgements:**

RosettaCommons is supported by NIH R01 GM073151 to Kuhlman, NSF, the Packard Foundation, the Beckman Foundation, the Alfred P. Sloan Foundation, and the Simons Foundation. This work was also supported by 100,000,000 CPU-hour donation from Google Inc to Conway; 125,760,000 CPU-hour allocation on the Mira and Theta supercomputers through the Innovative and Novel Computational Impact on Theory and Experiment (INCITE) program to Baker, DiMaio, Leaver-Fay, and Mulligan. This research used resources of the Argonne Leadership Computing Facility, which is a DOE Office of Science user facility supported under Contract DE-AC02-06CH11357. AHA 18POST34080422 to Kuenze; AMED J-PRIDE JP18fm0208022h to Kuroda; Biltema Foundation to Correia; Boehringer Ingelheim Fonds to Norn; Computing was performed using resources of the Argonne Leadership Computing Facility at Argonne National Laboratory which is supported by the Office of Science of the US to Conway; DFG KU 3510/1-1 to Kuenze; DP120100561 to Huber; DP150100383 to Huber, Pilla; Defense Threat Reduction Agency to King; EMBO long-term fellowship ALTF 698-2011 to Stein; EPFL-Fellows - H2020 Marie Skłodowska-Curie to Bonet; European Research Council Grant 310873 to Schueler-Furman, Alam; European Research Council Grant 310873 to Sedan, Marcu; European Research Council Starting grant - 716058 to Correia, Scheck; FT0991709 to Huber; Foundation of Knut and Alice Wallenberg 20160023 to Malmström; Hertz Foundation Fellowship to Alford; Howard Hughes Medical Institute to Baker; Hyak supercomputer system supported in part by the University of Washington eScience Institute to Baker and DiMaio labs; Israel Science Foundation 2017717 to Schueler-Furman, Alam; Japan Society for the Promotion of Science JP17K18113 to Kuroda; MCB1330760 to Khare; Marie Curie International Outgoing Fellowship FP7-PEOPLE-2011-IOF 298976 to Marcos; National Science Centre, Poland, 2018/29/B/ST6/01989 to Gront; NIAID T32AI007244 to Adolf-Bryfogle; NIAID U19 AI117905 to Sevy; NIEHS P42ES004699 to Siegel; NIGMS Ruth L Kirschstein National Research Service Award T32GM008268 to Conway; NIGMS T32 GM007628 to Bender; NIH 1R35 GM122579 to Das; NIH 1UH2CA203780 to Cooper, Khatib; NIH 5F32GM110899-02 to Linsky; NIH F31GM123616 to Jeliaskov; NIH F32CA189246 to Labonte; NIH P01 U19AI117905, R01 AI113867, UM1 AI100663 to Schief; NIH R00 GM120388 to Horowitz; NIH R01 AI143997 to Sgourakis; NIH R01 DK097376 to Meiler; NIH R01 GM073960 to Kuhlman; NIH R01 GM076324 to Siegel; NIH R01 GM078221 to Gray; NIH R01 GM080403 to Meiler; NIH R01 GM084453 to Dunbrack; NIH R01 GM088277 to Bradley; NIH R01 GM092802 to Baker; NIH R01 GM092802 to Baker; NIH R01 GM098101 to Kortemme; NIH R01 GM099842 to Meiler; NIH R01 GM099959 to Karanicolas; NIH R01 GM110089 to Kortemme; NIH R01 GM117189 to Kortemme; NIH R01 GM117968 to Kuhlman; NIH R01 GM121487 to Bradley; NIH R01 GM123089 to DiMaio; NIH R01 GM126299 to Pierce; NIH R01 GM127578 to Gray; NIH R01 GM099827 to Bystroff; NIH R01 GM092802 to Baker; NIH R01 HL122010 to Meiler; NIH R01067553 to Kuhlman; NIH R01 GM078221 to Gray; NIH R01084433 to Baker; NIH R01088277 to Thyme; NIH R21 AI121799 to Meiler; NIH R21 CA219847 to Das; NIH R21 GM102716 to Das; NIH R35 GM122517 to Dunbrack; NIH R35 GM125034 to Sgourakis; NIH RL1CA133832 to Baker; NIH U19 AI117905 to Meiler; NIH/NCI Cancer Center Support Grant P30 CA006927 to Karanicolas; NSF 1507736 to Gray; NSF 1627539 to Siegel; NSF 1629879 to Cooper; NSF 1805510 to Siegel; NSF 1827246 to Siegel; NSF CHE 1305874 to Meiler; NSF CHE 1750666 to Lindert; NSF CISE 1629811 to Meiler; NSF CNS-1629811 to Meiler; NSF DBI-1262182 to Kortemme; NSF DBI-1564692 to Kortemme; NSF DMR 1507736 to Gray; NSF GRF DGE-1433187 to Rubenstein; NSF Graduate Research Fellowship to Alford, Kappel, Koepnick, Thyme; NSF MCB1330760 to Khare; NSF MCB1716623 to Khare; Open Philanthropy to Coventry; PhRMA Informatics Pre-Doctoral Fellowship U22879-001 to Smith; PhRMA foundation Predoctoral Fellowship to Fu; RosettaCommons to Goldschmidt, Rubenstein, DiMaio, Cooper, Watkins, Szegedy, Geniesse, Blacklock, Das, Khare, Koehler Leman, Kappel; SEB is funded by a Career Award at the Scientific Interface from Burroughs Wellcome Fund to Boyken; Simons Foundation to Mulligan, Bonneau, Renfrew, Koehler Leman; Stanford Graduate Fellowship to Kappel; Starter Grant from the European Research Council to Lapidoth; Swiss National Science Foundation - NCCR Molecular Systems Engineering 51NF40-141825 to Correia; Swiss National Science Foundation 310030\_163139 to Correia; Swiss National Science Foundation SNF 200021 160188 to Malmström, Khakzad; UCSF/UCB Graduate Program in Bioengineering to Pan; USA-Israel Binational Science Foundation 2009418 to Raveh, Zimmerman, London; USA-Israel Binational Science Foundation 2009418, 2015207 to Schueler-Furman,



Alam; USA-Israel Binational Science Foundation 2015207 to Khrumushin; Washington Research Foundation Innovation Postdoctoral Fellowship to Weitzner; XSEDE, which is supported by NSF ACI-1548562.

**Author contributions:**

JKL wrote the manuscript with help from BDW. All authors edited and approved the manuscript and were substantially involved in developing the methods described, either by conception of the ideas or by implementing the methods into Rosetta. The idea for this paper was conceived by RB.

**Conflicts of Interest:**

Yifan Song, Indigo C. King, Steven M. Lewis, Brandon Frenz, and Ryan Pavlovicz are currently employed at Cyrus Biotechnology with granted stock options. Cyrus Biotechnology distributes the Rosetta software, which may include methods described in this paper. Baker, Malmström, Bonneau, Meiler, Kuhlman, Siegel, Gray, Gront, Karanicolas, Kortemme and Bradley are unpaid board members of the Rosetta Commons. Under institutional participation agreements between the University of Washington, acting on behalf of the Rosetta Commons, their respective institutions may be entitled to a portion of revenue received on licensing Rosetta software including programs described here. As members of the Scientific Advisory Board of Cyrus Biotechnology, Baker and Gray are granted stock options. Cyrus Biotechnology distributes the Rosetta software, which may include methods described in this paper. Brian D. Weitzner and Scott E. Boyken hold equity in Lyell Immunopharma. The content of this manuscript is relevant to work performed at Lyell.

## References:

1. Koehler Leman, J., Weitzner, B. D., ... & Bonneau, R. Better together: 20+ years of scientific software development in the Rosetta macromolecular modeling suite. *Prep.* (2019).
2. Leaver-Fay, A., Tyka, M., Lewis, S. M., Lange, O. F., Thompson, J. M., Jacak, R., Kaufman, K., Renfrew, P. D., Smith, C. A., Sheffler, W., Davis, I. W., Cooper, S., Treuille, A., Mandell, D. J., Richter, F., Ban, Y.-E. A., Fleishman, S. J., Corn, J. E., Kim, D. E., Berrondo, M., Mentzer, S., Popovic, Z., Havranek, J. J., Karanicolas, J., Das, R., Meiler, J., Kortemme, T., Gray, J. J., Kuhlman, B., Baker, D. & Bradley, P. ROSETTA3: An Object-Oriented Software Suite for the Simulation and Design of Macromolecules. *Methods Enzymol.* **487**, 545–74 (2011).
3. Alford, R. F., Leaver-Fay, A., Jeliaskov, J. R., O'Meara, M. J., Dimaio, F. P., Park, H., Shapovalov, M. V., Renfrew, P. D., Mulligan, V. K., Kappel, K., Labonte, J. W., Pacella, M. S., Bonneau, R., Bradley, P., Dunbrack, R. L., Das, R., Baker, D., Kuhlman, B., Kortemme, T. & Gray, J. J. The Rosetta all-atom energy function for macromolecular modeling and design. *J. Chem. Theory Comput.* **13**, 1–35 (2017).
4. Park, H., Bradley, P., Greisen, P., Liu, Y., Mulligan, V. K., Kim, D. E., Baker, D. & DiMaio, F. Simultaneous Optimization of Biomolecular Energy Functions on Features from Small Molecules and Macromolecules. *J. Chem. Theory Comput.* **12**, 6201–6212 (2016).
5. Chaudhury, S., Lyskov, S. & Gray, J. J. PyRosetta: a script-based interface for implementing molecular modeling algorithms using Rosetta. *Bioinformatics* **26**, 689–691 (2010).
6. Fleishman, S. J., Leaver-Fay, A., Corn, J. E., Strauch, E.-M. M., Khare, S. D., Koga, N., Ashworth, J., Murphy, P., Richter, F., Lemmon, G., Meiler, J. & Baker, D. RosettaScripts: A scripting language interface to the Rosetta Macromolecular modeling suite. *PLoS One* **6**, 1–10 (2011).
7. Cooper, S., Khatib, F., Treuille, A., Barbero, J., Lee, J., Beenen, M., Leaver-Fay, A., Baker, D., Popović, Z. & Players, F. Predicting protein structures with a multiplayer online game. *Nature* **466**, 756–760 (2010).
8. Bender, B. J., Cisneros, A., Duran, A. M., Finn, J. A., Fu, D., Lokits, A. D., Mueller, B. K., Sangha, A. K., Sauer, M. F., Sevy, A. M., Sliwoski, G., Sheehan, J. H., Dimaio, F., Meiler, J. & Moretti, R. Protocols for Molecular Modeling with Rosetta3 and RosettaScripts. *Biochemistry* [acs.biochem.6b00444](https://doi.org/10.1021/acs.biochem.6b00444) (2016). doi:10.1021/acs.biochem.6b00444
9. Gront, D., Kulp, D. W., Vernon, R. M., Strauss, C. E. M. & Baker, D. Generalized Fragment Picking in Rosetta: Design, Protocols and Applications. *PLoS One* **6**, e23294 (2011).
10. Song, Y., Dimaio, F., Wang, R. Y., Kim, D. E., Miles, C., Brunette, T. J., Thompson, J. & Baker, D. Supplemental Information High-Resolution Comparative Modeling with RosettaCM. **21**,
11. Kamisetty, H., Ovchinnikov, S. & Baker, D. Assessing the utility of coevolution-based residue-residue contact predictions in a sequence- and structure-rich era. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 15674–9 (2013).
12. Ovchinnikov, S., Park, H., Varghese, N., Huang, P.-S., Pavlopoulos, G. A., Kim, D. E., Kamisetty, H., Kypides, N. C. & Baker, D. Protein structure determination using metagenome sequence data. *Science (80-. )*. **355**, 294–298 (2017).
13. Park, H., Ovchinnikov, S., Kim, D. E., Dimaio, F. & Baker, D. Protein homology model refinement by large-scale energy optimization. *Proc. Natl. Acad. Sci. U. S. A.* **115**, 3054–3059 (2018).
14. Tyka, M. D., Keedy, D. A., André, I., Dimaio, F., Song, Y., Richardson, D. C., Richardson, J. S. & Baker, D. Alternate states of proteins revealed by detailed energy landscape mapping. *J. Mol. Biol.* **405**, 607–18 (2011).
15. Friedland, G. D., Linares, A. J., Smith, C. A. & Kortemme, T. A simple model of backbone flexibility improves modeling of side-chain conformational variability. *J. Mol. Biol.* **380**, 757–74 (2008).
16. Kapp, G. T., Liu, S., Stein, A., Wong, D. T., Remenyi, A., Yeh, B. J., Fraser, J. S., Taunton, J., Lim, W. A. & Kortemme, T. Control of protein signaling using a computationally designed GTPase/GEF orthogonal pair. *Proc. Natl. Acad. Sci.* **109**, 5277–5282 (2012).
17. Stein, A. & Kortemme, T. Improvements to robotics-inspired conformational sampling in rosetta. *PLoS One* **8**, e63090 (2013).
18. Pacella, M. S., Koo, D. C. E., Thottungal, R. A. & Gray, J. J. Using the RosettaSurface Algorithm to Predict Protein Structure at Mineral Surfaces. *Methods Enzymol.* **532**, 343–366 (2013).
19. Lubin, J. H., Pacella, M. S. & Gray, J. J. A Parametric Rosetta Energy Function Analysis with LK Peptides on SAM Surfaces. *Langmuir* **34**, 5279–5289 (2018).
20. Marze, N. A., Roy Burman, S. S., Sheffler, W. & Gray, J. J. Efficient flexible backbone protein–

- protein docking for challenging targets. *Bioinformatics* **34**, 3461–3469 (2018).
21. Roy Burman, S. S., Yovanno, R. A. & Gray, J. J. Flexible Backbone Assembly and Refinement of Symmetrical Homomeric Complexes. *Structure* **0**, (2019).
  22. DiMaio, F., Leaver-Fay, A., Bradley, P., Baker, D. & André, I. Modeling Symmetric Macromolecular Structures in Rosetta3. *PLoS One* **6**, e20450 (2011).
  23. Marcos, E., Basanta, B., Chidyausiku, T. M., Tang, Y., Oberdorfer, G., Liu, G., Swapna, G. V. T., Guan, R., Silva, D.-A., Dou, J., Pereira, J. H., Xiao, R., Sankaran, B., Zwart, P. H., Montelione, G. T. & Baker, D. Principles for designing proteins with cavities formed by curved  $\beta$  sheets. *Science* **355**, 201–206 (2017).
  24. Dang, B., Wu, H., Mulligan, V. K., Mravic, M., Wu, Y., Lemmin, T., Ford, A., Silva, D.-A., Baker, D. & DeGrado, W. F. De novo design of covalently constrained mesosize protein scaffolds with unique tertiary structures. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 10852–10857 (2017).
  25. Hauri, S., Khakzad, H., Happonen, L., Teleman, J., Malmström, J. & Malmström, L. Rapid determination of quaternary protein structures in complex biological samples. *Nat. Commun.* **10**, 192 (2019).
  26. Labonte, J. W., Adolf-Bryfogle, J., Schief, W. R. & Gray, J. J. Residue-centric modeling and design of saccharide and glycoconjugate structures. *J. Comput. Chem.* **38**, 276–287 (2017).
  27. Alam, N., Goldstein, O., Xia, B., Porter, K. A., Kozakov, D. & Schueler-Furman, O. High-resolution global peptide-protein docking using fragments-based PIPER-FlexPepDock. *PLoS Comput. Biol.* **13**, e1005905 (2017).
  28. Meiler, J. & Baker, D. RosettaLigand: protein-small molecule docking with full side-chain flexibility. *Proteins* **65**, 538–48 (2006).
  29. DeLuca, S., Khar, K. & Meiler, J. Fully Flexible Docking of Medium Sized Ligand Libraries with RosettaLigand. *PLoS One* **10**, e0132508 (2015).
  30. Davis, I. W. & Baker, D. RosettaLigand Docking with Full Ligand and Receptor Flexibility. *J. Mol. Biol.* **385**, 381–392 (2009).
  31. Fu, D. Y. & Meiler, J. Predictive Power of Different Types of Experimental Restraints in Small Molecule Docking: A Review. *J. Chem. Inf. Model.* **58**, 225–233 (2018).
  32. Fu, D. Y. & Meiler, J. RosettaLigandEnsemble: A Small-Molecule Ensemble-Driven Docking Approach. *ACS Omega* **3**, 3655–3664 (2018).
  33. Johnson, D. K. & Karanicolas, J. Selectivity by Small-Molecule Inhibitors of Protein Interactions Can Be Driven by Protein Surface Fluctuations. *PLOS Comput. Biol.* **11**, e1004081 (2015).
  34. Johnson, D. K. & Karanicolas, J. Druggable Protein Interaction Sites Are More Predisposed to Surface Pocket Formation than the Rest of the Protein Surface. *PLoS Comput. Biol.* **9**, e1002951 (2013).
  35. Johnson, D. K. & Karanicolas, J. Ultra-High-Throughput Structure-Based Virtual Screening for Small-Molecule Inhibitors of Protein–Protein Interactions. *J. Chem. Inf. Model.* **56**, 399–411 (2016).
  36. Gowthaman, R., Miller, S. A., Rogers, S., Khowsathit, J., Lan, L., Bai, N., Johnson, D. K., Liu, C., Xu, L., Anbanandam, A., Aubé, J., Roy, A. & Karanicolas, J. DARC: Mapping Surface Topography by Ray-Casting for Effective Virtual Screening at Protein Interaction Sites. *J. Med. Chem.* **59**, 4152–4170 (2016).
  37. Khar, K. R., Goldschmidt, L. & Karanicolas, J. Fast Docking on Graphics Processing Units via Ray-Casting. *PLoS One* **8**, e70661 (2013).
  38. Gowthaman, R., Lyskov, S. & Karanicolas, J. DARC 2.0: Improved Docking and Virtual Screening at Protein Interaction Sites. *PLoS One* **10**, e0131612 (2015).
  39. Jacobs, T. M., Williams, B., Williams, T., Xu, X., Eletsky, A., Federizon, J. F., Szyperski, T. & Kuhlman, B. Design of structurally distinct proteins using strategies inspired by evolution. **352**, (2016).
  40. Guffy, S. L., Teets, F. D., Langlois, M. I. & Kuhlman, B. Protocols for Requirement-Driven Protein Design in the Rosetta Modeling Program. *J. Chem. Inf. Model.* **58**, 895–901 (2018).
  41. Lapidoth, G., Khersonsky, O., Lipsh, R., Dym, O., Albeck, S., Rogotner, S. & Fleishman, S. J. Highly active enzymes by automated combinatorial backbone assembly and sequence design. *Nat. Commun.* **9**, 2780 (2018).
  42. Huang, P.-S., Ban, Y.-E. A., Richter, F., Andre, I., Vernon, R., Schief, W. R. & Baker, D. RosettaRemodel: A Generalized Framework for Flexible Backbone Protein Design. *PLoS One* **6**,

- e24109 (2011).
43. Blacklock, K. M., Yang, L., Mulligan, V. K. & Khare, S. D. A computational method for the design of nested proteins by loop-directed domain insertion. *Proteins Struct. Funct. Bioinforma.* **86**, 354–369 (2018).
  44. Bhardwaj, G., Mulligan, V. K., Bahl, C. D., Gilmore, J. M., Harvey, P. J., Cheneval, O., Buchko, G. W., Pulavarti, S. V. S. R. K., Kaas, Q., Eletsky, A., Huang, P.-S., Johnsen, W. A., Greisen, P. J., Rocklin, G. J., Song, Y., Linsky, T. W., Watkins, A., Rettie, S. A., Xu, X., Carter, L. P., Bonneau, R., Olson, J. M., Coutsiadis, E., Correnti, C. E., Szyperski, T., Craik, D. J. & Baker, D. Accurate de novo design of hyperstable constrained peptides. *Nature* **538**, 329–335 (2016).
  45. Leaver-Fay, A., Jacak, R., Stranges, P. B. & Kuhlman, B. A Generic Program for Multistate Protein Design. *PLoS One* **6**, e20937 (2011).
  46. Sevy, A. M., Jacobs, T. M., Crowe, J. E. & Meiler, J. Design of Protein Multi-specificity Using an Independent Sequence Search Reduces the Barrier to Low Energy Sequences. *PLoS Comput. Biol.* **11**, e1004300 (2015).
  47. Sevy, A. M., Wu, N. C., Gilchuk, I. M., Parrish, E. H., Burger, S., Yousif, D., Nagel, M. B. M., Schey, K. L., Wilson, I. A., Crowe, J. E. & Meiler, J. Multistate design of influenza antibodies improves affinity and breadth against seasonal viruses. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 1597–1602 (2019).
  48. Sauer, M. F., Sevy, A. M., Crowe, J. E. & Meiler, J. Manuscript submitted. (2019).
  49. Marcos, E. & Silva, D.-A. Essentials of *de novo* protein design: Methods and applications. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* **8**, e1374 (2018).
  50. Correia, B. E., Bates, J. T., Loomis, R. J., Baneyx, G., Carrico, C., Jardine, J. G., Rupert, P., Correnti, C., Kalyuzhnyi, O., Vittal, V., Connell, M. J., Stevens, E., Schroeter, A., Chen, M., MacPherson, S., Serra, A. M., Adachi, Y., Holmes, M. A., Li, Y., Klevit, R. E., Graham, B. S., Wyatt, R. T., Baker, D., Strong, R. K., Crowe, J. E., Johnson, P. R. & Schief, W. R. Proof of principle for epitope-focused vaccine design. *Nature* **507**, 201–206 (2014).
  51. Bonet, J., Wehrle, S., Schriever, K., Yang, C., Billet, A., Sesterhenn, F., Scheck, A., Sverrisson, F., Veselkova, B., Vollers, S., Lourman, R., Villard, M., Rosset, S., Krey, T. & Correia, B. E. Rosetta FunFolDes - A general framework for the computational design of functional proteins. *PLoS Comput. Biol.* **14**, e1006623 (2018).
  52. Kroncke, B. M., Duran, A. M., Mendenhall, J. L., Meiler, J., Blume, J. D. & Sanders, C. R. Documentation of an Imperative To Improve Methods for Predicting Membrane Protein Stability. *Biochemistry* **55**, 5002–5009 (2016).
  53. Kortemme, T. & Baker, D. A simple physical model for binding energy hot spots in protein-protein complexes. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 14116–21 (2002).
  54. Kortemme, T., Kim, D. E. & Baker, D. Computational alanine scanning of protein-protein interfaces. *Sci. STKE* **2004**, pl2 (2004).
  55. Ó Conchúir, S., Barlow, K. A., Pache, R. A., Ollikainen, N., Kundert, K., O'Meara, M. J., Smith, C. A. & Kortemme, T. A Web Resource for Standardized Benchmark Datasets, Metrics, and Rosetta Protocols for Macromolecular Modeling and Design. *PLoS One* **10**, e0130433 (2015).
  56. Barlow, K. A., Ó Conchúir, S., Thompson, S., Suresh, P., Lucas, J. E., Heinonen, M. & Kortemme, T. Flex ddG: Rosetta Ensemble-Based Estimation of Changes in Protein-Protein Binding Affinity upon Mutation. *J. Phys. Chem. B* **122**, 5389–5399 (2018).
  57. Smith, C. A. & Kortemme, T. Backrub-like backbone simulation recapitulates natural protein conformational variability and improves mutant side-chain prediction. *J. Mol. Biol.* **380**, 742–56 (2008).
  58. Dou, J., Doyle, L., Jr Greisen, P., Schena, A., Park, H., Johnsson, K., Stoddard, B. L. & Baker, D. Sampling and energy evaluation challenges in ligand binding protein design. *Protein Sci.* **26**, 2426–2437 (2017).
  59. Smith, C. A. & Kortemme, T. Structure-based prediction of the peptide sequence space recognized by natural and synthetic PDZ domains. *J. Mol. Biol.* **402**, 460–74 (2010).
  60. Smith, C. A. & Kortemme, T. Predicting the tolerated sequences for proteins and protein interfaces using RosettaBackrub flexible backbone design. *PLoS One* **6**, e20451 (2011).
  61. Ollikainen, N., de Jong, R. M. & Kortemme, T. Coupling Protein Side-Chain and Backbone Flexibility Improves the Re-design of Protein-Ligand Specificity. *PLoS Comput. Biol.* **11**, e1004335 (2015).

62. Crick, F. H. C. The Fourier transform of a coiled-coil. *Acta Crystallogr.* **6**, 685–689 (1953).
63. Lu, P., Min, D., DiMaio, F., Wei, K. Y., Vahey, M. D., Boyken, S. E., Chen, Z., Fallas, J. A., Ueda, G., Sheffler, W., Mulligan, V. K., Xu, W., Bowie, J. U. & Baker, D. Accurate computational design of multipass transmembrane proteins. *Science* (80-. ). **359**, 1042–1046 (2018).
64. Raveh, B., London, N. & Schueler-Furman, O. Sub-angstrom modeling of complexes between flexible peptides and globular proteins. *Proteins* **78**, 2029–40 (2010).
65. Raveh, B., London, N., Zimmerman, L. & Schueler-Furman, O. Rosetta FlexPepDock ab-initio: Simultaneous Folding, Docking and Refinement of Peptides onto Their Receptors. *PLoS One* **6**, e18934 (2011).
66. Kozakov, D., Brenke, R., Comeau, S. R. & Vajda, S. PIPER: an FFT-based protein docking program with pairwise potentials. *Proteins* **65**, 392–406 (2006).
67. London, N., Raveh, B., Movshovitz-Attias, D. & Schueler-Furman, O. Can self-inhibitory peptides be derived from the interfaces of globular protein-protein interactions? *Proteins* **78**, 3140–9 (2010).
68. Watkins, A. M. & Arora, P. S. Anatomy of  $\beta$ -strands at protein-protein interfaces. *ACS Chem. Biol.* **9**, 1747–54 (2014).
69. Jochim, A. L. & Arora, P. S. Systematic analysis of helical protein interfaces reveals targets for synthetic inhibitors. *ACS Chem. Biol.* **5**, 919–23 (2010).
70. Gavenonis, J., Sheneman, B. A., Siegert, T. R., Eshelman, M. R. & Kritzer, J. A. Comprehensive analysis of loops at protein-protein interfaces for macrocycle design. *Nat. Chem. Biol.* **10**, 716–22 (2014).
71. Sedan, Y., Marcu, O., Lyskov, S. & Schueler-Furman, O. Peptiderive server: derive peptide inhibitors from protein–protein interactions. *Nucleic Acids Res.* gkw385 (2016). doi:10.1093/nar/gkw385
72. Moretti, R., Lyskov, S., Das, R., Meiler, J. & Gray, J. J. Web-accessible molecular modeling with Rosetta: The Rosetta Online Server that Includes Everyone (ROSIE). *Protein Sci.* **27**, 259–268 (2018).
73. Hosseinzadeh, P., Bhardwaj, G., Mulligan, V. K., Shortridge, M. D., Craven, T. W., Pardo-Avila, F., Rettie, S. A., Kim, D. E., Silva, D.-A., Ibrahim, Y. M., Webb, I. K., Cort, J. R., Adkins, J. N., Varani, G. & Baker, D. Comprehensive computational design of ordered peptide macrocycles. *Science* (80-. ). **358**, 1461–1466 (2017).
74. Rubenstein, A. B., Pethe, M. A. & Khare, S. D. MFPred: Rapid and accurate prediction of protein-peptide recognition multispecificity using self-consistent mean field theory. *PLOS Comput. Biol.* **13**, e1005614 (2017).
75. Rohl, C. A., Strauss, C. E. M., Chivian, D. & Baker, D. Modeling structurally variable regions in homologous proteins with rosetta. *Proteins* **55**, 656–77 (2004).
76. Wang, C., Bradley, P. & Baker, D. Protein-Protein Docking with Backbone Flexibility. *J. Mol. Biol.* **373**, 503–519 (2007).
77. Canutescu, A. A. & Dunbrack, R. L. Cyclic coordinate descent: A robotics algorithm for protein loop closure. *Protein Sci.* **12**, 963–72 (2003).
78. Mandell, D. J., Coutsias, E. A. & Kortemme, T. Sub-angstrom accuracy in protein loop reconstruction by robotics-inspired conformational sampling. *Nat. Methods* **6**, 551–2 (2009).
79. Mandell, D. J. & Kortemme, T. Backbone flexibility in computational protein design. *Curr. Opin. Biotechnol.* **20**, 420–8 (2009).
80. Pan, X. & Kortemme, T. Manuscript in preparation. (2019).
81. Tyka, M. D., Jung, K. & Baker, D. Efficient sampling of protein conformational space using fast loop building and batch minimization on highly parallel computers. *J. Comput. Chem.* **33**, 2483–2491 (2012).
82. Marcos, E., Chidyausiku, T. M., McShan, A. C., Evangelidis, T., Nerli, S., Carter, L., Nivón, L. G., Davis, A., Oberdorfer, G., Tripsianes, K., Sgourakis, N. G. & Baker, D. De novo design of a non-local  $\beta$ -sheet protein with high stability and accuracy. *Nat. Struct. Mol. Biol.* **25**, 1028–1034 (2018).
83. Sircar, A., Kim, E. T. & Gray, J. J. RosettaAntibody: antibody variable region homology modeling server. *Nucleic Acids Res.* **37**, W474–479 (2009).
84. Weitzner, B. D., Kuroda, D., Marze, N., Xu, J. & Gray, J. J. Blind prediction performance of RosettaAntibody 3.0: Grafting, relaxation, kinematic loop modeling, and full CDR optimization. *Proteins Struct. Funct. Bioinforma.* **82**, 1611–1623 (2014).
85. Weitzner, B. D., Jeliakov, J. R., Lyskov, S., Marze, N., Kuroda, D., Frick, R., Adolf-Bryfogle, J.,



- Biswas, N., Dunbrack, R. L. & Gray, J. J. Modeling and docking of antibody structures with Rosetta. *Nat. Protoc.* **12**, 401–416 (2017).
86. Sivasubramanian, A., Sircar, A., Chaudhury, S. & Gray, J. J. Toward high-resolution homology modeling of antibody F<sub>v</sub> regions and application to antibody-antigen docking. *Proteins Struct. Funct. Bioinforma.* **74**, 497–514 (2009).
87. Marze, N. A., Lyskov, S. & Gray, J. J. Improved prediction of antibody V<sub>L</sub>–V<sub>H</sub> orientation. *Protein Eng. Des. Sel.* **29**, 409–418 (2016).
88. Finn, J. A., Koehler Leman, J., Willis, J. R., Cisneros, A., Crowe, J. E. & Meiler, J. Improving Loop Modeling of the Antibody Complementarity-Determining Region 3 Using Knowledge-Based Restraints. *PLoS One* **11**, e0154811 (2016).
89. Weitzner, B. D. & Gray, J. J. Accurate Structure Prediction of CDR H3 Loops Enabled by a Novel Structure-Based C-Terminal Constraint. *J. Immunol.* **198**, 505–515 (2017).
90. Sircar, A., Sanni, K. A., Shi, J. & Gray, J. J. Analysis and modeling of the variable region of camelid single-domain antibodies. *J. Immunol.* **186**, 6357–67 (2011).
91. DeKosky, B. J., Lungu, O. I., Park, D., Johnson, E. L., Charab, W., Chrysostomou, C., Kuroda, D., Ellington, A. D., Ippolito, G. C., Gray, J. J. & Georgiou, G. Large-scale sequence and structural comparisons of human naive and antigen-experienced antibody repertoires. *Proc. Natl. Acad. Sci.* **113**, E2636–E2645 (2016).
92. Jeliakov, J. R., Sljoka, A., Kuroda, D., Tsuchimura, N., Katoh, N., Tsumoto, K. & Gray, J. J. Repertoire Analysis of Antibody CDR-H3 Loops Suggests Affinity Maturation Does Not Typically Result in Rigidification. *Front. Immunol.* **9**, 413 (2018).
93. Norn, C. H., Lapidoth, G. & Fleishman, S. J. High-accuracy modeling of antibody structures by a search for minimum-energy recombination of backbone fragments. *Proteins* **85**, 30–38 (2017).
94. Lapidoth, G., Parker, J., Prilusky, J. & Fleishman, S. J. AbPredict 2: a server for accurate and unstrained structure prediction of antibody variable domains. *Bioinformatics* (2018). doi:10.1093/bioinformatics/bty822
95. Nerli, S. & Sgourakis, N. G. Manuscript in preparation. (2019).
96. Toor, J. S., Rao, A. A., McShan, A. C., Yarmarkovich, M., Nerli, S., Yamaguchi, K., Madejska, A. A., Nguyen, S., Tripathi, S., Maris, J. M., Salama, S. R., Haussler, D. & Sgourakis, N. G. A Recurrent Mutation in Anaplastic Lymphoma Kinase with Distinct Neopeptide Conformations. *Front. Immunol.* **9**, 99 (2018).
97. Gowthaman, R. & Pierce, B. G. TCRmodel: high resolution modeling of T cell receptors from sequence. *Nucleic Acids Res.* **46**, W396–W401 (2018).
98. Tran, E., Robbins, P. F., Lu, Y.-C., Prickett, T. D., Gartner, J. J., Jia, L., Pasetto, A., Zheng, Z., Ray, S., Groh, E. M., Kriley, I. R. & Rosenberg, S. A. T-Cell Transfer Therapy Targeting Mutant KRAS in Cancer. *N. Engl. J. Med.* **375**, 2255–2262 (2016).
99. Pierce, B. G. & Weng, Z. A flexible docking approach for prediction of T cell receptor-peptide-MHC complexes. *Protein Sci.* **22**, 35–46 (2013).
100. Pierce, B. G., Vreven, T. & Weng, Z. Modeling T cell receptor recognition of CD1-lipid and MR1-metabolite complexes. *BMC Bioinformatics* **15**, 319 (2014).
101. Sircar, A. & Gray, J. J. SnugDock: paratope structural optimization during antibody-antigen docking compensates for errors in antibody homology models. *PLoS Comput. Biol.* **6**, e1000644 (2010).
102. Adolf-Bryfogle, J., Kalyuzhnyi, O., Kubitz, M., Weitzner, B. D., Hu, X., Adachi, Y., Schief, W. R. & Dunbrack, R. L. RosettaAntibodyDesign (RABD): A general framework for computational antibody design. *PLOS Comput. Biol.* **14**, e1006112 (2018).
103. North, B., Lehmann, A. & Dunbrack, R. L. A New Clustering of Antibody CDR Loop Conformations. *J. Mol. Biol.* **406**, 228–256 (2011).
104. King, C., Garza, E. N., Mazor, R., Linehan, J. L., Pastan, I., Pepper, M. & Baker, D. Removing T-cell epitopes with computational protein design. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 8577–82 (2014).
105. Nivón, L. G., Bjelic, S., King, C. & Baker, D. Automating human intuition for protein design. *Proteins* **82**, 858–66 (2014).
106. Lapidoth, G. D., Baran, D., Pszolla, G. M., Norn, C., Alon, A., Tyka, M. D. & Fleishman, S. J. AbDesign: An algorithm for combinatorial backbone design guided by natural conformations and sequences. *Proteins* **83**, 1385–406 (2015).

107. Baran, D., Pszolla, M. G., Lapidoth, G. D., Norn, C., Dym, O., Unger, T., Albeck, S., Tyka, M. D. & Fleishman, S. J. Principles for computational design of binding antibodies. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 10900–10905 (2017).
108. Wang, R. Y.-R., Kudryashev, M., Li, X., Egelman, E. H., Basler, M., Cheng, Y., Baker, D. & DiMaio, F. De novo protein structure determination from near-atomic-resolution cryo-EM maps. *Nat. Methods* **12**, 335–8 (2015).
109. Frenz, B., Walls, A. C., Egelman, E. H., Veessler, D. & DiMaio, F. RosettaES: a sampling strategy enabling automated interpretation of difficult cryo-EM maps. *Nat. Methods* **14**, 797–800 (2017).
110. DiMaio, F., Echols, N., Headd, J. J., Terwilliger, T. C., Adams, P. D. & Baker, D. Improved low-resolution crystallographic refinement with Phenix and Rosetta. *Nat. Methods* **10**, 1102–4 (2013).
111. DiMaio, F., Song, Y., Li, X., Brunner, M. J., Xu, C., Conticello, V., Egelman, E., Marlovits, T. C., Cheng, Y. & Baker, D. Atomic-accuracy models from 4.5-Å cryo-electron microscopy data with density-guided iterative local refinement. *Nat. Methods* **12**, 361–5 (2015).
112. Wang, R. Y.-R., Song, Y., Barad, B. A., Cheng, Y., Fraser, J. S. & DiMaio, F. Automated structure refinement of macromolecular assemblies from cryo-EM maps using Rosetta. *Elife* **5**, (2016).
113. Frenz, B., Rämisch, S., Borst, A. J., Walls, A. C., Adolf-Bryfogle, J., Schief, W. R., Veessler, D. & DiMaio, F. Automatically Fixing Errors in Glycoprotein Structures with Rosetta. *Structure* **0**, (2018).
114. Nerli, S. & Sgourakis, N. G. CS-ROSETTA. *Methods Enzymol.* (2018). doi:10.1016/BS.MIE.2018.07.005
115. Rohl, C. A. & Baker, D. De novo determination of protein backbone structure from residual dipolar couplings using Rosetta. *J. Am. Chem. Soc.* **124**, 2723–9 (2002).
116. Schmitz, C., Vernon, R., Otting, G., Baker, D. & Huber, T. Protein Structure Determination from Pseudocontact Shifts Using ROSETTA. *J. Mol. Biol.* **416**, 668–677 (2012).
117. Yagi, H., Pilla, K. B., Maleckis, A., Graham, B., Huber, T. & Otting, G. Three-dimensional protein fold determination from backbone amide pseudocontact shifts generated by lanthanide tags at multiple sites. *Structure* **21**, 883–890 (2013).
118. Pilla, K. B., Otting, G. & Huber, T. Pseudocontact Shift-Driven Iterative Resampling for 3D Structure Determinations of Large Proteins. *J. Mol. Biol.* **428**, 522–532 (2016).
119. Lange, O. F. & Baker, D. Resolution-adapted recombination of structural features significantly improves sampling in restraint-guided structure calculation. *Proteins Struct. Funct. Bioinforma.* **80**, 884–895 (2012).
120. Bowers, P. M., Strauss, C. E. M. & Baker, D. De novo protein structure determination using sparse NMR data. 311–318 (2000).
121. Meiler, J. & Baker, D. Rapid protein fold determination using unassigned NMR data. *Proc. Natl. Acad. Sci. U. S. A.* **100**, 15404–9 (2003).
122. Raman, S., Raman, S., Lange, O. F., Rossi, P., Tyka, M., Wang, X., Aramini, J., Liu, G., Ramelot, T. A., Eletsky, A., Szyperski, T., Kennedy, M. A., Prestegard, J., Montelione, G. T. & Baker, D. NMR Structure Determination for Larger Proteins Using Backbone-Only Data. **1014**, (2010).
123. Lange, O. F., Rossi, P., Sgourakis, N. G., Song, Y., Lee, H.-W., Aramini, J. M., Ertekin, a., Xiao, R., Acton, T. B., Montelione, G. T. & Baker, D. Determination of solution structures of proteins up to 40 kDa using CS-Rosetta with sparse NMR data from deuterated samples. *Proc. Natl. Acad. Sci.* **109**, 10873–10878 (2012).
124. Reichel, K., Fiset, O., Braun, T., Lange, O. F., Hummer, G. & Schäfer, L. V. Systematic evaluation of CS-Rosetta for membrane protein structure prediction with sparse NOE restraints. *Proteins* **85**, 812–826 (2017).
125. Sgourakis, N. G., Lange, O. F., DiMaio, F., André, I., Fitzkee, N. C., Rossi, P., Montelione, G. T., Bax, A. & Baker, D. Determination of the structures of symmetric protein oligomers from NMR chemical shifts and residual dipolar couplings. *J. Am. Chem. Soc.* **133**, 6288–98 (2011).
126. Rossi, P., Shi, L., Liu, G., Barbieri, C. M., Lee, H. W., Grant, T. D., Luft, J. R., Xiao, R., Acton, T. B., Snell, E. H., Montelione, G. T., Baker, D., Lange, O. F. & Sgourakis, N. G. A hybrid NMR/SAXS-based approach for discriminating oligomeric protein interfaces using Rosetta. *Proteins Struct. Funct. Bioinforma.* **83**, 309–317 (2015).
127. Demers, J.-P., Habenstein, B., Loquet, A., Kumar Vasa, S., Giller, K., Becker, S., Baker, D., Lange, A. & Sgourakis, N. G. High-resolution structure of the Shigella type-III secretion needle by solid-state NMR and cryo-electron microscopy. *Nat. Commun.* **5**, 4976 (2014).
128. Thompson, J. M., Sgourakis, N. G., Liu, G., Rossi, P., Tang, Y., Mills, J. L., Szyperski, T.,

- Montelione, G. T. & Baker, D. Accurate protein structure modeling using sparse NMR data and homologous structure information. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 9875–9880 (2012).
129. Braun, T., Koehler Leman, J. & Lange, O. F. Combining Evolutionary Information and an Iterative Sampling Strategy for Accurate Protein Structure Prediction. *PLoS Comput. Biol.* **11**, (2015).
  130. Evangelidis, T., Nerli, S., Nováček, J., Brereton, A. E., Karplus, P. A., Dotas, R. R., Venditti, V., Sgourakis, N. G. & Tripsianes, K. Automated NMR resonance assignments and structure determination using a minimal set of 4D spectra. *Nat. Commun.* **9**, 384 (2018).
  131. Lange, O. F. Automatic NOESY assignment in CS-RASREC-Rosetta. *J. Biomol. NMR* **59**, 147–159 (2014).
  132. Kuenze, G., Bonneau, R., Leman, J. K. & Meiler, J. Integrative protein modeling in RosettaNMR from sparse paramagnetic restraints. *bioRxiv* 597872 (2019). doi:10.1101/597872
  133. Aprahamian, M. L., Chea, E. E., Jones, L. M. & Lindert, S. Rosetta Protein Structure Prediction from Hydroxyl Radical Protein Footprinting Mass Spectrometry Data. *Anal. Chem.* **90**, 7721–7729 (2018).
  134. Aprahamian, M. L. & Lindert, S. Utility of Covalent Labeling Mass Spectrometry Data in Protein Structure Prediction with Rosetta. *J. Chem. Theory Comput.* **19**, 1021/acs.jctc.9b00101 (2019). doi:10.1021/acs.jctc.9b00101
  135. Watkins, A. M., Geniesse, C., Kladwang, W., Zakrevsky, P., Jaeger, L. & Das, R. Blind prediction of noncanonical RNA structure at atomic accuracy. *Sci. Adv.* **4**, eaar5316 (2018).
  136. Sripakdeevong, P., Kladwang, W. & Das, R. An enumerative stepwise ansatz enables atomic-accuracy RNA loop modeling. *Proc. Natl. Acad. Sci.* **108**, 20573–20578 (2011).
  137. Das, R. Atomic-Accuracy Prediction of Protein Loop Structures through an RNA-Inspired Ansatz. *PLoS One* **8**, e74830 (2013).
  138. Chou, F.-C., Sripakdeevong, P., Dibrov, S. M., Hermann, T. & Das, R. Correcting pervasive errors in RNA crystallography through enumerative structure prediction. *Nat. Methods* **10**, 74–76 (2013).
  139. Chou, F.-C., Echols, N., Terwilliger, T. C. & Das, R. in 269–282 (Humana Press, New York, NY, 2016). doi:10.1007/978-1-4939-2763-0\_17
  140. Das, R., Karanicolas, J. & Baker, D. Atomic accuracy in predicting and designing noncanonical RNA structure. *Nat. Methods* **7**, 291–294 (2010).
  141. Cheng, C. Y., Chou, F.-C. & Das, R. Modeling Complex RNA Tertiary Folds with Rosetta. *Methods Enzymol.* **553**, 35–64 (2015).
  142. Sripakdeevong, P., Cevec, M., Chang, A. T., Erat, M. C., Ziegeler, M., Zhao, Q., Fox, G. E., Gao, X., Kennedy, S. D., Kierzek, R., Nikonowicz, E. P., Schwalbe, H., Sigel, R. K. O., Turner, D. H. & Das, R. Structure determination of noncanonical RNA motifs guided by <sup>1</sup>H NMR chemical shifts. *Nat. Methods* **11**, 413–416 (2014).
  143. Kappel, K. & Das, R. Sampling Native-like Structures of RNA-Protein Complexes through Rosetta Folding and Docking. *Structure* **27**, 140-151.e5 (2019).
  144. Kappel, K., Jarmoskaite, I., Vaidyanathan, P. P., Greenleaf, W. J., Herschlag, D. & Das, R. Blind tests of RNA–protein binding affinity prediction. *Proc. Natl. Acad. Sci.* **116**, 8336–8341 (2019).
  145. Kappel, K., Liu, S., Larsen, K. P., Skiniotis, G., Puglisi, E. V., Puglisi, J. D., Zhou, Z. H., Zhao, R. & Das, R. De novo computational RNA modeling into cryo-EM maps of large ribonucleoprotein complexes. *Nat. Methods* **15**, 947–954 (2018).
  146. Thyme, S. B., Jarjour, J., Takeuchi, R., Havranek, J. J., Ashworth, J., Scharenberg, A. M., Stoddard, B. L. & Baker, D. Exploitation of binding energy for catalysis and design. *Nature* **461**, 1300–1304 (2009).
  147. Ashworth, J., Havranek, J. J., Duarte, C. M., Sussman, D., Monnat, R. J., Stoddard, B. L. & Baker, D. Computational redesign of endonuclease DNA binding and cleavage specificity. *Nature* **441**, 656–659 (2006).
  148. Ashworth, J., Taylor, G. K., Havranek, J. J., Quadri, S. A., Stoddard, B. L. & Baker, D. Computational reprogramming of homing endonuclease specificity at multiple adjacent base pairs. *Nucleic Acids Res.* **38**, 5601–5608 (2010).
  149. Havranek, J. J. & Harbury, P. B. Automated design of specificity in molecular recognition. *Nat. Struct. Biol.* **10**, 45–52 (2003).
  150. Thyme, S. B., Boissel, S. J. S., Arshiya Quadri, S., Nolan, T., Baker, D. A., Park, R. U., Kusak, L., Ashworth, J. & Baker, D. Reprogramming homing endonuclease specificity through computational design and directed evolution. *Nucleic Acids Res.* **42**, 2564–2576 (2014).

151. Thyme, S. B., Baker, D. & Bradley, P. Improved Modeling of Side-Chain–Base Interactions and Plasticity in Protein–DNA Interface Design. *J. Mol. Biol.* **419**, 255–274 (2012).
152. Yanover, C. & Bradley, P. Extensive protein and DNA backbone sampling improves structure-based specificity prediction for C2H2 zinc fingers. *Nucleic Acids Res.* **39**, 4564–76 (2011).
153. Ashworth, J. & Baker, D. Assessment of the optimization of affinity and specificity at protein–DNA interfaces. *Nucleic Acids Res.* **37**, e73 (2009).
154. Song, Y., Dimaio, F., Wang, R. Y.-R. R., Kim, D. E., Miles, C., Brunette, T., Thompson, J. & Baker, D. High-resolution comparative modeling with RosettaCM. *Structure* **21**, 1735–1742 (2013).
155. Thyme, S. B., Song, Y., Brunette, T. J., Szeto, M. D., Kusak, L., Bradley, P. & Baker, D. Massively parallel determination and modeling of endonuclease substrate specificity. *Nucleic Acids Res.* **42**, 13839–13852 (2014).
156. Overington, J. P., Al-Lazikani, B. & Hopkins, A. L. How many drug targets are there? *Nat. Rev. Drug Discov.* **5**, 993–6 (2006).
157. Koehler Leman, J., Ulmschneider, M. B. & Gray, J. J. Computational modeling of membrane proteins. *Proteins Struct. Funct. Bioinforma.* **83**, 1–24 (2015).
158. Yarov-Yarovoy, V., Schonbrun, J. & Baker, D. Multipass membrane protein structure prediction using Rosetta. *Proteins* **62**, 1010–1025 (2006).
159. Yarov-Yarovoy, V., Decaen, P. G., Westenbroek, R. E., Pan, C.-Y. Y. C.-Y., Scheuer, T., Baker, D. & Catterall, W. a. Structural basis for gating charge movement in the voltage sensor of a sodium channel. *Proc. Natl. Acad. Sci.* **109**, E93–E102 (2012).
160. Barth, P., Schonbrun, J. & Baker, D. Toward high-resolution prediction and design of transmembrane helical protein structures. **2007**, (2007).
161. Alford, R. F., Koehler Leman, J., Weitzner, B. D., Duran, A. M., Tilley, D. C., Elazar, A. & Gray, J. J. An Integrated Framework Advancing Membrane Protein Modeling and Design. *PLoS Comput. Biol.* **11**, e1004398 (2015).
162. Baugh, E. H., Lyskov, S., Weitzner, B. D. & Gray, J. J. Real-time PyMOL visualization for Rosetta and PyRosetta. *PLoS One* **6**, e21931 (2011).
163. Koehler Leman, J., Mueller, B. K. & Gray, J. J. Expanding the toolkit for membrane protein modeling in Rosetta. *Bioinformatics* **11**, 1–3 (2016).
164. Koehler Leman, J., Lyskov, S. & Bonneau, R. Computing structure-based lipid accessibility of membrane proteins with mp\_lipid\_acc in RosettaMP. *BMC Bioinformatics* **18**, 115 (2017).
165. Koehler Leman, J. & Bonneau, R. A novel domain assembly routine for creating full-length models of membrane proteins from known domain structures. *Biochemistry* acs.biochem.7b00995 (2017). doi:10.1021/acs.biochem.7b00995
166. Varki, A. Biological roles of oligosaccharides: all of the theories are correct. *Glycobiology* **3**, 97–130 (1993).
167. Varki, A., Cummings, R. D., Esko, J. D., Freeze, H. H., Stanley, P., Bertozzi, C. R., Hart, G. W. & Etzler, M. E. *Essentials of Glycobiology. Essentials Glycobiol.* (Cold Spring Harbor Laboratory Press, 2009).
168. Nivedha, A. K., Thieker, D. F., Makeneni, S., Hu, H. & Woods, R. J. Vina-Carb: Improving Glycosidic Angles during Carbohydrate Docking. *J. Chem. Theory Comput.* **12**, 892–901 (2016).
169. Leaver-Fay, A., O'Meara, M. J., Tyka, M., Jacak, R., Song, Y., Kellogg, E. H., Thompson, J., Davis, I. W., Pache, R. A., Lyskov, S., Gray, J. J., Kortemme, T., Richardson, J. S., Havranek, J. J., Snoeyink, J., Baker, D. & Kuhlman, B. Scientific benchmarks for guiding macromolecular energy function improvement. *Methods Enzymol.* **523**, 109–43 (2013).
170. Jorgensen, W. L., Jorgensen, W. L., Maxwell, D. S. & Tirado-rives, J. Development and testing of the OPLS all-atom force field on conformational energetics and properties of organic liquids. *J. AM. CHEM. SOC* **11225–11236** (1996). at <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.334.2959>
171. Radzicka, A. & Wolfenden, R. Comparing the polarities of the amino acids: side-chain distribution coefficients between the vapor phase, cyclohexane, 1-octanol, and neutral aqueous solution. *Biochemistry* **27**, 1664–1670 (1988).
172. O'Meara, M. J., Leaver-Fay, A., Tyka, M. D., Stein, A., Houlihan, K., DiMaio, F., Bradley, P., Kortemme, T., Baker, D., Snoeyink, J. & Kuhlman, B. Combined Covalent-Electrostatic Model of Hydrogen Bonding Improves Structure Prediction with Rosetta. *J. Chem. Theory Comput.* **11**, 609–622 (2015).



173. Conway, P., Tyka, M. D., DiMaio, F., Konerding, D. E. & Baker, D. Relaxation of backbone bond geometry improves protein energy landscape modeling. *Protein Sci.* **23**, 47–55 (2014).
174. Park, H., Lee, H. & Seok, C. High-resolution protein-protein docking by global optimization: recent advances and future challenges. *Curr. Opin. Struct. Biol.* **35**, 24–31 (2015).
175. Kellogg, E. H., Leaver-Fay, A. & Baker, D. Role of conformational sampling in computing mutation-induced changes in protein structure and stability. *Proteins Struct. Funct. Bioinforma.* **79**, 830–838 (2011).
176. Mills, J. H., Khare, S. D., Bolduc, J. M., Forouhar, F., Mulligan, V. K., Lew, S., Seetharaman, J., Tong, L., Stoddard, B. L. & Baker, D. Computational Design of an Unnatural Amino Acid Dependent Metalloprotein with Atomic Level Accuracy. *J. Am. Chem. Soc.* **135**, 13393–13399 (2013).
177. Mulligan, V. K. Manuscript in preparation. (2019).
178. Leaver-Fay, A., Butterfoss, G. L., Snoeyink, J. & Kuhlman, B. Maintaining solvent accessible surface area under rotamer substitution for protein design. *J. Comput. Chem.* **28**, 1336–41 (2007).
179. Rubenstein, A. B., Blacklock, K., Nguyen, H., Case, D. A. & Khare, S. D. Systematic Comparison of Amber and Rosetta Energy Functions for Protein Structure Evaluation. *J. Chem. Theory Comput.* **14**, 6015–6025 (2018).
180. Boyken, S. E., Chen, Z., Groves, B., Langan, R. A., Oberdorfer, G., Ford, A., Gilmore, J. M., Xu, C., DiMaio, F., Pereira, J. H., Sankaran, B., Seelig, G., Zwart, P. H. & Baker, D. De novo design of protein homo-oligomers with modular hydrogen-bond network-mediated specificity. *Science* **352**, 680–7 (2016).
181. Chen, Z., Boyken, S. E., Jia, M., Busch, F., Flores-Solis, D., Bick, M. J., Lu, P., VanAernum, Z. L., Sahasrabudhe, A., Langan, R. A., Bermeo, S., Brunette, T. J., Mulligan, V. K., Carter, L. P., DiMaio, F., Sgourakis, N. G., Wysocki, V. H. & Baker, D. Programmable design of orthogonal protein heterodimers. *Nature* **565**, 106–111 (2019).
182. Maguire, J. B., Boyken, S. E., Baker, D. & Kuhlman, B. Rapid Sampling of Hydrogen Bond Networks for Computational Protein Design. *J. Chem. Theory Comput.* **14**, 2751–2760 (2018).
183. Pavlovicz, R. E., Park, H. & DiMaio, F. Efficient consideration of coordinated water molecules improves computational protein-protein and protein-ligand docking. *bioRxiv* 618603 (2019). doi:10.1101/618603
184. Gray, J. J., Chaudhury, S., Lyskov, S., and Labonte, J. W. The PyRosetta Interactive Platform for Protein Structure Prediction and Design: A Set of Educational Modules. (2014). at <<http://www.amazon.com/PyRosetta-Interactive-Platform-Structure-Prediction/dp/1500968277>>
185. Schenkelberg, C. D. & Bystroff, C. InteractiveROSETTA: A graphical user interface for the PyRosetta protein modeling suite. *Bioinformatics* (2015). doi:10.1093/bioinformatics/btv492
186. Kleffner, R., Flatten, J., Leaver-Fay, A., Baker, D., Siegel, J. B., Khatib, F. & Cooper, S. Foldit Standalone: a video game-derived protein structure manipulation interface using Rosetta. *Bioinformatics* **33**, 2765–2767 (2017).
187. Cooper, S., Sterling, A. L. R., Kleffner, R., Silversmith, W. M. & Siegel, J. B. Repurposing citizen science games as software tools for professional scientists. in *Proc. 13th Int. Conf. Found. Digit. Games - FDG '18* 1–6 (ACM Press, 2018). doi:10.1145/3235765.3235770
188. Khatib, F., Cooper, S., Tyka, M. D., Xu, K., Makedon, I., Popovic, Z., Baker, D. & Players, F. Algorithm discovery by protein folding game players. *Proc. Natl. Acad. Sci. U. S. A.* **108**, 18949–53 (2011).
189. Dsilva, L., Mittal, S., Koepnick, B., Flatten, J., Cooper, S. & Horowitz, S. Creating custom Foldit puzzles for teaching biochemistry. *Biochem. Mol. Biol. Educ.* **47**, 133–139 (2019).
190. Lyskov, S., Chou, F.-C., Conchúir, S. Ó., Der, B. S., Drew, K., Kuroda, D., Xu, J., Weitzner, B. D., Renfrew, P. D., Sripakdeevong, P., Borgo, B., Havranek, J. J., Kuhlman, B., Kortemme, T., Bonneau, R., Gray, J. J. & Das, R. Serverification of molecular modeling applications: the Rosetta Online Server that Includes Everyone (ROSIE). *PLoS One* **8**, e63906 (2013).