

Supporting Information Appendix for: “Blind tests of RNA nearest neighbor energy prediction”

Fang-Chieh Chou¹, Wipapat Kladwang¹, Kalli Kappel², and Rhiju Das^{1,2,3*}

¹ Department of Biochemistry, Stanford University, Stanford, CA 94305, USA

² Biophysics Program, Stanford University, Stanford, CA 94305, USA

³ Department of Physics, Stanford University, Stanford, CA 94305, USA

* Correspondence to: Rhiju Das. Phone: (650) 723-5976. Fax: (650) 723-6783. E-mail: rhiju@stanford.edu.

| | |
|--|----|
| Supporting Methods..... | 2 |
| Nearest neighbor Parameter Estimation with RECCES..... | 2 |
| Nearest neighbor Parameters with Symmetric G-C pairs..... | 2 |
| Nearest Neighbor Parameter Estimation, General Case..... | 4 |
| The Dangling-End Nearest neighbor Parameters..... | 5 |
| Free Energy Computation and Score Function Reweighting..... | 6 |
| Details of Simulated Tempering Monte Carlo simulations..... | 11 |
| Details of Reweighting and Training Scheme..... | 13 |
| Electrostatic interactions across stacked bases..... | 15 |
| Simple single-conformation methods..... | 17 |
| Hydrogen-bond Counting Model..... | 17 |
| Single-conformation Rosetta Scoring Method..... | 19 |
| Optical Melting Measurements on Helices with D-U Base pairs..... | 20 |
| Supporting Results..... | 21 |
| Accuracies of terminal base pair penalty calculations..... | 21 |
| Supporting Tables..... | 22 |
| Supporting Figures..... | 27 |
| References for Supporting Information Appendix..... | 29 |

Supporting Methods

Nearest neighbor Parameter Estimation with RECCES

In this section, we describe computation of NN parameters under the RECCES framework. First, we review how NN parameters can be determined from either experimental measurements or computational estimates of the folding free energies of individual helix motifs. Hereafter, ‘folding free energies’ refers to the free energy of association of two separate strands into a folded A-form helix at 1 M standard state. First, we discuss a simple example involving association of two GC segments with each other, in the context of a longer helix. Second, we describe the calculations of NN parameters for the general case where the helix strand segments have different sequences as well as calculations for dangling-ends, again in terms of well-defined folding free energies of specific complexes. The presented relations apply to both experimental and computational studies. Third, for computational approaches, we show how each of these folding free energies can be computed based on the free energies of the complex and of the separate strands (in a ‘random coil’ state). Last, we describe the RECCES sampling framework, which uses simulated-tempering Monte Carlo to efficiently sample the density of states of and then evaluate the free energy of each single-strand and helix molecule.

Nearest neighbor Parameters with Symmetric G-C pairs

Here we compute a NN parameter of the stacked pair $\begin{pmatrix} 5'GC \\ 3'CG \end{pmatrix}$ using the folding free energies of two helix motifs with different lengths. The NN model assumes that the folding free energy of an RNA helix can be decomposed into the sum of the contribution of each NN fragment. For example,

$$\begin{aligned}\Delta G_f \left(\begin{matrix} 5'GG \\ 3'CC \end{matrix} \right) &= \Delta G_{NN} \left(\begin{matrix} 5'GG \\ 3'CC \end{matrix} \right) + \Delta G_{init} \\ \Delta G_f \left(\begin{matrix} 5'GGC \\ 3'CCG \end{matrix} \right) &= \Delta G_{NN} \left(\begin{matrix} 5'GG \\ 3'CC \end{matrix} \right) + \Delta G_{NN} \left(\begin{matrix} 5'GC \\ 3'CG \end{matrix} \right) + \Delta G_{init}\end{aligned}\tag{1}$$

where ΔG_f is the folding free energy of the given helix; the terms ΔG_{NN} are the NN parameters; and ΔG_{init} accounts for the entropic penalty of initiating the helix with the first G-C base pair. Hereafter, the free energies are computed at standard states in which each molecule is at 1 M concentration and the temperature is 37 °C, conditions at which most values of NN parameters are tabulated in the experimental literature. (Similar expressions for the temperature dependence of energetics lead to relations involving ΔH and ΔS . Calculating these parameters from simulations requires understanding the dependence of solvation and other physical effects on temperature and will not be considered herein.)

In the above example, the folding free energy expression of the second helix contains the $\Delta G_{NN} \left(\begin{matrix} 5'GC \\ 3'CG \end{matrix} \right)$ parameter, which we wish to determine from experimental measurements or computer simulations.

Taking the difference of the two folding free energies above, we have:

$$\Delta G_{NN} \left(\begin{matrix} 5'GC \\ 3'CG \end{matrix} \right) = \Delta G_f \left(\begin{matrix} 5'GGC \\ 3'CCG \end{matrix} \right) - \Delta G_f \left(\begin{matrix} 5'GG \\ 3'CC \end{matrix} \right)\tag{2}$$

Note that ΔG_{init} cancels out in the above equation. Calculating ΔG_{init} requires accounting for the translational and rotational entropy lost upon helix association and requires separate computations beyond the scope of the present work.

Nearest Neighbor Parameter Estimation, General Case

The example above illustrates the evaluation of NN parameters involving two helix-associating segments with the same two-nucleotide sequence. The more general case involves taking into account an additional parameter, the terminal penalty, as discussed next.

As an example, we describe the computation of the nearest neighbor parameter $\Delta G_{NN} \left(\begin{smallmatrix} 5'GA3' \\ 3'CU5' \end{smallmatrix} \right)$.

Following the procedure in eqs. (1)-(2) above, this can be evaluated from folding free energies

$\Delta G_f \left(\begin{smallmatrix} 5'GGA \\ 3'CCU \end{smallmatrix} \right)$ and $\Delta G_f \left(\begin{smallmatrix} 5'GG \\ 3'CC \end{smallmatrix} \right)$. The relationships assumed by the NN model are:

$$\begin{aligned} \Delta G_f \left(\begin{smallmatrix} 5'GGA \\ 3'CCU \end{smallmatrix} \right) &= \Delta G_{NN} \left(\begin{smallmatrix} 5'GG \\ 3'CC \end{smallmatrix} \right) + \Delta G_{NN} \left(\begin{smallmatrix} 5'GA \\ 3'CU \end{smallmatrix} \right) + \Delta G_{init} + \Delta G_{Terminal-AU} \\ \Delta G_f \left(\begin{smallmatrix} 5'GG \\ 3'CC \end{smallmatrix} \right) &= \Delta G_{NN} \left(\begin{smallmatrix} 5'GG \\ 3'CC \end{smallmatrix} \right) + \Delta G_{init} \end{aligned} \quad (3)$$

Here $\Delta G_{Terminal-AU}$ is the terminal contribution for having an A-U pair at the terminal of the helix instead of G-C. Determining this additional term requires additional folding free energies, as follows. The

definition of the NN model requires $\Delta G_{NN} \left(\begin{smallmatrix} 5'GA3' \\ 3'CU5' \end{smallmatrix} \right)$ to be the same as $\Delta G_{NN} \left(\begin{smallmatrix} 5'UC3' \\ 3'AG5' \end{smallmatrix} \right)$, flipping the upper

and lower sequence. Analogous to eq. (3), we can write the parameter in terms of $\Delta G_f \left(\begin{smallmatrix} 5'GUC \\ 3'CAG \end{smallmatrix} \right)$ and

$\Delta G_f \left(\begin{smallmatrix} 5'GU \\ 3'CA \end{smallmatrix} \right)$:

$$\begin{aligned}
\Delta G_f \left(\begin{array}{c} 5'GUC \\ 3'CAG \end{array} \right) &= \Delta G_{NN} \left(\begin{array}{c} 5'GU \\ 3'CA \end{array} \right) + \Delta G_{NN} \left(\begin{array}{c} 5'UC \\ 3'AG \end{array} \right) + \Delta G_{init} \\
\Delta G_f \left(\begin{array}{c} 5'GU \\ 3'CA \end{array} \right) &= \Delta G_{NN} \left(\begin{array}{c} 5'GU \\ 3'CA \end{array} \right) + \Delta G_{init} + \Delta G_{Terminal-AU}
\end{aligned} \tag{4}$$

The NN parameter and the terminal A-U contribution can now be computed in terms of measurable folding free energies ΔG_f of helix association as:

$$\begin{aligned}
\Delta G_{NN} \left(\begin{array}{c} 5'GA \\ 3'CU \end{array} \right) &= \frac{1}{2} \left[\Delta G_f \left(\begin{array}{c} 5'GGA \\ 3'CCU \end{array} \right) - \Delta G_f \left(\begin{array}{c} 5'GG \\ 3'CC \end{array} \right) + \Delta G_f \left(\begin{array}{c} 5'GUC \\ 3'CAG \end{array} \right) - \Delta G_f \left(\begin{array}{c} 5'GU \\ 3'CA \end{array} \right) \right] \\
\Delta G_{Terminal-AU} &= \frac{1}{2} \left[\Delta G_f \left(\begin{array}{c} 5'GGA \\ 3'CCU \end{array} \right) - \Delta G_f \left(\begin{array}{c} 5'GG \\ 3'CC \end{array} \right) - \Delta G_f \left(\begin{array}{c} 5'GUC \\ 3'CAG \end{array} \right) + \Delta G_f \left(\begin{array}{c} 5'GU \\ 3'CA \end{array} \right) \right]
\end{aligned} \tag{5}$$

Other NN parameters and terminal contributions can be similarly evaluated in terms of folding free energies of two-base-pair and three-base-pair helices. In our calculations, we checked that using helices of different lengths or with a different first base pair led to systematic errors of 0.26 kcal/mol or less in final NN values (Table S1).

The Dangling-End Nearest neighbor Parameters

The other NN parameters considered herein are for single-nucleotide dangling ends. These dangling-end parameters contribute to the folding free energies of complexes such as $\Delta G_f \left(\begin{array}{c} 5'GAA \\ 3'CU \end{array} \right)$:

$$\begin{aligned}
\Delta G_f \left(\begin{array}{c} 5'GAA \\ 3'CU \end{array} \right) &= \Delta G_{NN} \left(\begin{array}{c} 5'GA \\ 3'CU \end{array} \right) + \Delta G_{NN} \left(\begin{array}{c} 5'AA \\ 3'U \end{array} \right) + \Delta G_{init} + \Delta G_{Terminal-AU} \\
\Delta G_f \left(\begin{array}{c} 5'GA \\ 3'CU \end{array} \right) &= \Delta G_{NN} \left(\begin{array}{c} 5'GA \\ 3'CU \end{array} \right) + \Delta G_{init} + \Delta G_{Terminal-AU}
\end{aligned} \tag{6}$$

Similar to above [eq. (2) and (5)], the dangling-end NN parameter $\Delta G_{NN} \left(\begin{smallmatrix} 5'AA \\ 3'U \end{smallmatrix} \right)$ can be estimated from the difference of folding free energies for complexes with and without the dangling-ends:

$$\Delta G_{NN} \left(\begin{smallmatrix} 5'AA \\ 3'U \end{smallmatrix} \right) = \Delta G_f \left(\begin{smallmatrix} 5'GAA \\ 3'CU \end{smallmatrix} \right) - \Delta G_f \left(\begin{smallmatrix} 5'GA \\ 3'CU \end{smallmatrix} \right) \quad (7)$$

As a brief note of clarification, the experimental measurements of the dangling-end parameters were presented in a 1995 study, before the A-U terminal contribution was introduced into the NN model (1); eq. (6) is the appropriate update to include the terminal contribution. In either case, however, the simple expression (7) applies to determine the dangling end NN parameters from folding free energies, and the tabulated computational values below use this shared expression.

Free Energy Computation and Score Function Reweighting

In the above sections, we expressed the NN parameters as linear combinations of the folding free energies of strands folding into helices, which can be measured experimentally or estimated through computation. In our computational approach, the folding free energy of a helix is defined as the difference in the free energies of the helix and of the two separated single-strand molecules, i.e. the two states involved in the folding equilibrium (**Error! Reference source not found.a**):

$$\begin{aligned}
& \text{GGC} + \text{GCC} \rightleftharpoons \begin{matrix} 5'\text{GGC} \\ 3'\text{CCG} \end{matrix} \\
\Delta G_f \left(\begin{matrix} 5'\text{GGC} \\ 3'\text{CCG} \end{matrix} \right) &= G \left(\begin{matrix} 5'\text{GGC} \\ 3'\text{CCG} \end{matrix} \right) - G(\text{GGC}) - G(\text{GCC}) \\
& \text{GG} + \text{CC} \rightleftharpoons \begin{matrix} 5'\text{GG} \\ 3'\text{CC} \end{matrix} \\
\Delta G_f \left(\begin{matrix} 5'\text{GG} \\ 3'\text{CC} \end{matrix} \right) &= G \left(\begin{matrix} 5'\text{GG} \\ 3'\text{CC} \end{matrix} \right) - G(\text{GG}) - G(\text{CC})
\end{aligned} \tag{8}$$

Here the terms $G \left(\begin{matrix} 5'\text{GGC} \\ 3'\text{CCG} \end{matrix} \right)$, $G(\text{GGC})$, etc. are not free energy differences (ΔG) but are instead the 'raw' free energies of the helices and single-strands, each defined by an integral over the system's conformational space:

$$\begin{aligned}
G &= -k_B T \ln Z \\
Z &= \iint \dots \int d\theta_1 d\theta_2 \dots d\theta_n \exp \left(-\frac{E(\theta_1, \theta_2, \dots, \theta_n)}{k_B T} \right)
\end{aligned} \tag{9}$$

Here the partition function Z is the integration of the Boltzmann factor over all accessible torsion angles θ_i of the molecule, E is the internal energy of the molecule at a conformation specified by the torsion angles, k_B is the Boltzmann constant, and T is the system's temperature. In the folded state, because the energy of a single molecule is independent of translational and rotational degrees of freedom, those degrees of freedom only contribute to a constant factor to the free energy of the system, which is omitted in the above expressions. For bimolecular systems (e.g., the two disassociated strands), the system free energy does depend on the relative positions and orientations of the molecules, but the translational/rotational entropy from these terms cancel out in the determination of the NN parameters [see, e.g., eqs. (2) and (7) above] considered herein.

While eq. (9) can theoretically be used directly to compute the molecule's free energy, it is challenging to integrate the expression over all torsion angles, especially when the system has a large number of degrees of freedom. Instead we compute the partition function by estimating the molecule's density of states (DOS), $g(E)$:

$$\begin{aligned}
 Z &= V_p \int_{-\infty}^{\infty} dE g(E) \exp\left(-\frac{E}{kT}\right) \\
 g(x) &= \frac{\int \dots \int d\theta_1 \dots d\theta_n \delta(x - E(\theta_1, \dots, \theta_n))}{V_p} \\
 V_p &= \int \dots \int d\theta_1 \dots d\theta_n
 \end{aligned} \tag{10}$$

Here δ is the Dirac delta function; V_p is the total phase space volume available to the molecule, which can be calculated exactly (see below). Due to normalization by the V_p factor, the DOS $g(E)$ integrates to unity. The DOS describes the probability distribution of the molecule's energy at infinite temperature. Once determined, it allows the calculation of the partition function and hence free energies at any temperature. However, in practice, an infinite temperature simulation gives negligible sampling of the DOS for low energy states. Therefore, to estimate the DOS precisely at all temperatures, we carry out simulations of the conformational ensemble at different temperatures, providing estimates of $g(E)$ within different, overlapping temperature ranges, up to a different scaling factor for each simulation. Overlaying these distributions in overlapping energy regimes defines the unknown scaling factors at each temperature and yields a portrait of $g(E)$ across all temperature ranges, with the final overall normalization set by the property that $g(E)$ integrates to unity.

To carry out simulations of conformational ensembles, we used Metropolis Monte Carlo sampling, with states defined as follows. For a single RNA strand, we allowed all backbone and side-chain torsions to freely sample the entire range of $2\pi = 360^\circ$, except for the sugar puckers, which were only allowed to sample two conformations, the ideal 2'-endo and 3'-endo puckers (**Error! Reference source not found.**a, left panels). The phase space volume is $V_p = 2^n (2\pi)^n (2\pi)^{5(n-1)}$, where n is the number of nucleotides in the strand. Here 2^n represents the two sugar pucker forms, $(2\pi)^n$ represents the side-chain torsions (χ angle), and $(2\pi)^{5(n-1)}$ represent the backbone torsions connecting the sugars [five torsions ($\epsilon, \zeta, \alpha, \beta, \gamma$) for each connection, $n - 1$ connections between n nucleotides]. The 5' and 3' terminal phosphates were omitted in the simulated constructs. When included, these phosphates did not make stable interactions and instead gave constant entropic contributions that canceled out during the folding free energy evaluation.

For the helix state, we froze the relative position between the two strands by forcing the first base pair to take an ideal geometry (**Error! Reference source not found.**a, right panels). This constraint eliminates translation and rotational entropy contributions to computed helix free energy, which cancel out during the calculations of NN parameters; see, e.g., eq. (2). This cancellation allows the sampling to focus on estimating energy fluctuations and conformational entropy important for the considered NN parameters.

The sugar puckers of all nucleotides were restricted to the 3'-endo form, and the other torsions were allowed to sample values between $\pm 60^\circ$ around ideal A-form torsion angles ($\alpha = -64^\circ, \beta = 176^\circ, \gamma = 53^\circ, \epsilon = -150^\circ, \zeta = -71^\circ; \chi = 79^\circ$). The phase space volume over both strands is therefore $V_p = \left(\frac{2\pi}{3}\right)^{2n} \left(\frac{2\pi}{3}\right)^{2 \times 5(n-1)}$,

where n is the length of each strand in the helix. We note that the constraints described above provide our working definition of the helix state. We confirmed in separate calculations that changing the backbone torsion constraint from $\pm 60^\circ$ to $\pm 40^\circ$ and using alternative ideal pucker conformations, led to

negligible changes in computed free energies (0.08 and 0.13 kcal/mol error; Table S1). For both helix and strand states, we allowed the torsion angle for 2'-OH to sample 360°, leading to additional phase space volume contributions that canceled out during evaluation of the folding free energy.

For all Monte Carlo runs, we used a new application ('recces_turner') in Rosetta, with the computed energy based on the Rosetta scoring function. (See Supporting Results for command-lines and version numbers; Rosetta is freely available for academic users at www.rosettacommons.org.) The function is a linear combination of multiple component terms, including a Lennard-Jones potential, hydrogen bonding terms, an orientation-dependent solvation term, and long-range electrostatics (2). As noted above, we performed Monte Carlo simulations at multiple temperatures, ranging from room temperature to infinity. We applied the simulated tempering method (3) to facilitate the conformational search and barrier crossing during the Monte Carlo runs by allowing the temperature to vary during the simulation in a manner satisfying detailed balance (**Error! Reference source not found.**; see Supporting Methods for details). For each molecule, a typical simulation took less than an hour on a single CPU to generate up to 10 million Monte Carlo samples. The simulated tempering parameters were determined using short single-temperature simulations, as described by Huang et al. (4) We combined all simulations into one DOS using the Weighted Histogram Analysis Method (WHAM) (5, 6).

The comprehensive conformational ensemble obtained from the above simulation scheme allowed rapid tests of alternative energy functions based on different weights on the Rosetta score terms. Caching the contribution of each score term for each sampled conformation allowed rapid calculation of the entire simulation ensemble with different score-term weights; see also refs (7, 8). For each molecular ensemble, this reweighting step was a linear matrix operation that took less than 0.1 s on a

single CPU. **Error! Reference source not found.**c-d illustrates this procedure. This rapid reweighting enabled optimization of the weight set, by minimizing the difference between the computed NN parameters and the experimental values in a training set using a standard quasi-Newton minimization algorithm. Because the target cost function is not convex and has a large number of local minima, we repeated the minimization using thousands of randomly initialized starting weights. These separate minimization runs gave a collection of score functions, all compatible with the experimental NN parameters in the training set, and allowed an estimate of the systematic errors for the prediction of new NN parameters.

Details of Simulated Tempering Monte Carlo simulations

The free energy of each sequence (single-strand or helix) was evaluated using simulated tempering Monte Carlo (MC). For each sequence, we performed simulations at 7 temperatures T , with $k_B T = 0.8, 1, 1.4, 1.8, 3, 7$ and 30 Rosetta units (RU), where k_B is the Boltzmann constant. In later reweighting stages we calibrated the score so that one RU equals one $k_B T$ for $T = 310.15$ K (37 °C, at which most nearest neighbor parameters are tabulated). Before the simulated tempering run, we first performed short regular MC simulations (300,000 steps) at each of the seven temperatures, to determine parameters that govern switching between temperatures during the simulated tempering run. The simulated tempering parameters were computed by numerically solving the weight difference for each pair of neighbor temperatures such that the mean probability of moving upward and downward in the temperature ladder was the same (4). The initial values of these parameters, used by the numerical equation solver, were computed using the average energy for each temperature (9). With the computed simulated tempering parameters, we performed a long simulated tempering simulation (9,000,000 steps). In addition to simulated tempering, we also performed a regular MC simulation at infinite

temperature, to calibrate the full DOS profile. For all regular MC and simulated tempering simulations, acceptance rates for the conformational moves and temperature switches were above 10%. The simulations stored the total score and each separate score term component (Table S2, Table S5 and Table S6) for all sampled conformations, so that scores could be recomputed for different weight sets without rerunning the simulations. Example Rosetta command lines are given below:

1. Short pre-runs:

```
recces_turner -score:weights stepwise/rna/turner -n_cycle 300000 -seq1 gaa -temps 0.8 -  
out_prefix prerun
```

The above command performs a short single-temperature pre-run for the determination of the simulated tempering parameters. Here $k_B T = 0.8$ RU, and the simulated sequence is GAA (single strand).

2. Simulated tempering:

```
recces_turner -score:weights stepwise/rna/turner -seq1 gu -seq2 aac -n_cycle 9000000 -  
temps 0.8 1 1.4 1.8 3 7 30 -st_weights 0 8.04 15.39 17.78 17.76 14.96 11.81 -out_prefix  
ST
```

This command performs simulated tempering simulation on $\begin{matrix} 5'GU \\ 3'CAA \end{matrix}$ (dangling end).

3. Infinite temperature:

```
recces_turner -score:weights stepwise/rna/turner -seq1 gu -seq2 aac -n_cycle 300000 -  
temps -1 -out_prefix kT_inf
```

The “-1” temperature stands for infinite temperature, where all MC moves are accepted.

These simulation results were then combined into a single DOS using WHAM. First the simulation data at each temperature were aggregated into histograms with bin size of 0.1 RU. For the infinite temperature

simulation, the simulated conformation could have very high energy (> 1000 RU). Since these high energy conformations are not sampled at room temperature, and are important only for normalizing the density of states, we binned all conformations with scores higher than 800 RU into a single bin. We verified that using different bin sizes and score cutoffs in this procedure gave negligible changes to the results (<0.1 kcal/mol, Table S1). These histograms were then combined into one DOS using WHAM. During the combination calculation, WHAM assigned a weight to each energy bin (so all conformations in the same bin share the same weight). We recorded these weights for all conformations, which were needed for reweighting, as discussed next.

Details of Reweighting and Training Scheme

After collecting the conformational samples for all relevant sequences, we reweighted the Rosetta score function to minimize the prediction error on a limited training set of canonical and dangling end NN parameters. In this section we present the technical details of these reweighting and training steps.

For each sampled conformation, the Rosetta score is the weighted linear combination of all individual score terms:

$$Score(w_1, L, w_m) = \sum_{i=1}^m w_i s_i \quad (11)$$

Table S4 gives a short description of each score term; all were introduced in ref. (10) except for an electrostatic interaction between nucleobases (*stack_elec*), described below. Here w_j is the weight for each score term, and s_j is the value of the score term for the conformation. To reweight the score with a new set of weights, we updated the value of w_j in equation (11) to obtain a new score for each

conformation. The reweighted scores were then combined into DOS, using the WHAM weights computed previously (see section above). This reweighted DOS was then used to compute the free energy under the new weight set.

To optimize weight sets from training data (i.e. canonical and dangling-end parameters), we minimized the training error with respect to the score term weights. The training error was a weighted square error between the experimental and predicted NN parameters:

$$Error = \sum (G_{\text{canonical,expt}} - G_{\text{canonical,pred}})^2 + 0.1 \times \sum (G_{\text{dangling,expt}} - G_{\text{dangling,pred}})^2 + (G_{\text{terminal,expt}} - G_{\text{terminal,pred}})^2 \quad (12)$$

Each dangling end data point was given a weight of 0.1 compared to canonical data points, to avoid overfitting to the dangling end values, which were measured less accurately with fewer experimental measurements. During the training step, we minimized the error function (12) with respect to the score term weights, using the truncated Newton minimizer in Scipy (method *TNC* in `scipy.optimize.minimize`). To ensure the minimized weights were reasonable, we constrained the score term weights to be within certain ranges during minimization. Here the weights for *fa_atr* and *fa_rep* were constrained to be in [0.1, 20], weights for *hbond_sc* and *rna_torsion* were constrained to be in [0.5, 20], and the weight for *fa_intra_rep* was constrained to be in [0, 0.01]. All other score terms were constrained to be in [0, 20]. Before the minimization, the score term weights were randomly initialized to be between 1/10 to 10 times of the initial weights (Table S3). Minimizing with all the sampled data was computationally slow due to the large number of conformations being reweighted in each stage. Instead, we used only 1% of the data (randomly selected) for each minimization run. The final training errors for the minimized weight sets were computed using the full dataset. Because the target error function is not convex, the

minimization generated different results when initialized with different values. Rather than attempt an aggressive global optimization, we chose to assemble a large collection of locally optimized weight sets to help bracket systematic errors in predicting new energetic parameters. Here we repeated the minimization 9,544 times to obtain a diverse set of minimized weights (see Table S5 for example weights). These minimized weight sets were not all distinct; clustering analysis with a cluster radius of 0.4 (based on the Euclidean distance between the weight vectors) on 95% of the lowest error weight sets led to 1315 unique clusters. All the reweighting and minimization stages were carried out in separate Python scripts, also available in Rosetta (in subdirectory tools/rna_tools/recces).

Electrostatic interactions across stacked bases

Electrostatic interactions are critical components of the energetics of biomolecules. In previous Rosetta RNA modeling work, hydrogen bonding interactions were modeled using a potential derived from database hydrogen bond geometries and tested with quantum mechanical calculations (11).

Electrostatic contributions beyond this hydrogen bonding term were modeled in a limited manner (10, 12-14), through a weak carbon hydrogen bond potential (not included here due to lack of constraints from NN data; not shown) and a highly screened electrostatic repulsion between the backbone phosphates (*fa_elec_rna_phos_phos*)(10). Electrostatic interactions between bonded atoms were included implicitly in the RNA torsional potential, which was derived in a knowledge-based fashion using crystal structures in Protein Data Bank (PDB). The rather limited modeling of electrostatics followed the philosophy of Rosetta protein modeling, in which complex physical effects were left out until strong evidence from structure prediction or design tests have supported and helped calibrate the inclusion of these terms (15).

Recent studies of RNA model systems have suggested that the electrostatic interaction between the stacking bases might play an important role in the NN interactions of RNA (16). To test this hypothesis, we implemented a new score term called *stack_elec*. This new term models the electrostatic interaction between atom pairs in different bases through Coulomb's law with a distance-proportional dielectric (the same as the Rosetta *fa_elec* term), but the interaction is suppressed to zero for atom pairs whose inter-atom vector r is perpendicular to the base normal of the first or second base (parametrized by angles κ_1 and κ_2): $\frac{q_1 q_2}{r^2} (\cos^2 \kappa_1 + \cos^2 \kappa_2)$, where q_1 and q_2 are the atom partial charges. As a result of the orientation dependence, the electrostatic interactions between stacked but not co-planar hydrogen-bonded bases are captured in this term, and so *stack_elec* could be optimized separately from the knowledge-based Rosetta hydrogen bond potential. Without the orientation-dependent suppression (i.e., use of the original Rosetta *fa_elec* term), we observed strong repulsion between fixed positive charges of hydrogen atoms in, e.g., A-U pairs that could not be reconciled with the stability of the Watson-Crick arrangement. RECCES modeling of canonical stacked pairs with *fa_elec* instead of *stack_elec* gave worse RMSDs to canonical pair NN parameters (0.48 kcal/mol compared to 0.33 kcal/mol); and in simulations with both terms, optimization of these terms' weights recovered *stack_elec* at positive weight and *fa_elec* at zero weight. The partial charges for non-natural RNA residues were derived using the MATCH method (multipurpose atom-typer for CHARMM) (17). For canonical RNA residues, we used the default atomic partial charges in Rosetta for the *stack_elec* calculations. Calculations for NN parameters for canonical stacked pairs using MATCH charges instead of default Rosetta charges gave similar results, with an RMSD of 0.21 kcal/mol (comparable or less than other sources of systematic error; see Table S1).

Simple single-conformation methods

As baseline comparisons for the RECCES calculations above, we sought NN parameter predictions from phenomenological models that required negligible computation: a hydrogen-bond counting model and a single-conformation Rosetta scoring model. These two models used the same pipeline as RECCES to predict the NN parameters, by first computing the folding free energies for each helix motif, then combining these free energies into NN parameters. However instead of performing comprehensive ensemble sampling, these two models used simple approximations to evaluate the folding free energies.

Hydrogen-bond Counting Model

For the hydrogen-bond counting model, the folding free energy for a helix construct is

$$\Delta G_f(S_1, S_2) = N_{HB}(S_1, S_2)\Delta G_{HB} + N_{Dangle}(S_1, S_2)\Delta G_{Dangle} + [L(S_1) + L(S_2) - 2]\Delta G_{Stack} + [L(S_1) + L(S_2)]\Delta G_{Base} + \Delta G_{init} \quad (13)$$

Here S_1 and S_2 are the two single-strands in a helix construct. $N_{HB}(S_1, S_2)$ and $N_{Dangle}(S_1, S_2)$ are the number of hydrogen-bonds and dangling-ends in a fully associated helix conformation. For example,

$\begin{matrix} 5'GA \\ 3'CU \end{matrix}$ has 5 hydrogen bonds (3 for the G-C, plus 2 for the A-U) and no dangling-ends; $\begin{matrix} 5'GAA \\ 3'CU \end{matrix}$ also has 5

hydrogen bonds, plus 1 dangling-end. $L(S_1)$ and $L(S_2)$ are the lengths of the single-strand S_1 and S_2 .

ΔG_{HB} is the free energy contribution per hydrogen-bond, ΔG_{Dangle} is the free energy contribution for the conformational entropy of each dangling-end, ΔG_{Stack} is the free energy contribution for each base-base stacking ($L - 1$ base-base stacking per strand), and ΔG_{Base} is the entropy loss during folding for each base in the single-strands. This model assumes the folding free energy of a helix is a simple linear function of

the number of hydrogen bonds, as hydrogen-bonding is the dominant interaction in helix formation; the contributions of other relevant physical factors, such as base stacking and entropy loss upon helix formation, are approximated to be sequence independent.

We may further rewrite eq. (13) as follows:

$$\Delta G_f(S_1, S_2) = N_{HB}(S_1, S_2)\Delta G_{HB} + N_{Dangle}(S_1, S_2)\Delta G_{Dangle} + [L(S_1) + L(S_2)]\Delta G_S + \Delta G_{const} \quad (14)$$

Here $\Delta G_S = \Delta G_{Stack} + \Delta G_{Base}$ and $\Delta G_{const} = \Delta G_{init} - 2G_{Stack}$. In eq. (14), ΔG_{HB} , ΔG_{Dangle} , and ΔG_S are model parameters; ΔG_{const} is a constant factor that cancels out when we compute the NN parameters.

With the above folding free energy expressions, we can linearly combine them into NN parameters, as

we did in the RECESS framework. For example, in this model, the NN parameters for $\begin{matrix} 5'GC \\ 3'CG \end{matrix}$ and $\begin{matrix} 5'AA \\ 3'U \end{matrix}$

can be expressed as

$$\begin{aligned} \Delta G_{NN,HB} \left(\begin{matrix} 5'GC \\ 3'CG \end{matrix} \right) &= 3\Delta G_{HB} + 2\Delta G_S \\ \Delta G_{NN,HB} \left(\begin{matrix} 5'AA \\ 3'U \end{matrix} \right) &= \Delta G_{Dangle} + \Delta G_S \end{aligned} \quad (15)$$

The model parameters were determined by least-squares regression against the experimental NN parameters for canonical base pairs and dangling-ends (see Supporting Methods). The optimized parameters were $\Delta G_{HB} = -1.89$ kcal/mol, $\Delta G_S = +1.31$ kcal/mol, and $\Delta G_{Dangle} = -1.83$ kcal/mol.

Single-conformation Rosetta Scoring Method

Instead of simply counting the hydrogen bonds, the single-conformation Rosetta scoring method uses the Rosetta score to approximate the free energy of a construct. Similar to eq. (13), the model is

$$\Delta G_f(S_1, S_2) = k_{Rosetta} \text{Score}(S_1, S_2) + N_{Dangle}(S_1, S_2) \Delta G_{Dangle} + [L(S_1) + L(S_2)] \Delta G_s + \Delta G_{const} \quad (16)$$

Here, $\text{Score}(S_1, S_2)$ is the Rosetta score of a representative conformation of the target helix. This representative conformation was obtained by minimizing the helix conformation under the Rosetta score function. We used the standard Rosetta score function for RNA structure prediction (*rna_hires*).

$k_{Rosetta}$ is a model parameter that sets the scale between the Rosetta score and the folding free energy in kcal/mol; see also ref. (10).

For example, with this model, the NN parameters for $\begin{smallmatrix} 5'GC \\ 3'CG \end{smallmatrix}$ and $\begin{smallmatrix} 5'AA \\ 3'U \end{smallmatrix}$ can be expressed as

$$\begin{aligned} \Delta G_{NN} \left(\begin{smallmatrix} 5'GC \\ 3'CG \end{smallmatrix} \right) &= k_{Rosetta} \left[\text{Score} \left(\begin{smallmatrix} 5'GGC \\ 3'CCG \end{smallmatrix} \right) - \text{Score} \left(\begin{smallmatrix} 5'GG \\ 3'CC \end{smallmatrix} \right) \right] + 2\Delta G_s \\ \Delta G_{NN} \left(\begin{smallmatrix} 5'AA \\ 3'U \end{smallmatrix} \right) &= k_{Rosetta} \left[\text{Score} \left(\begin{smallmatrix} 5'GAA \\ 3'CU \end{smallmatrix} \right) - \text{Score} \left(\begin{smallmatrix} 5'GA \\ 3'CU \end{smallmatrix} \right) \right] + \Delta G_{Dangle} + \Delta G_s \end{aligned} \quad (17)$$

The corresponding model parameters, determined by training against experimental NN parameters with canonical base pairs and dangling-ends, are $k_{Rosetta} = 0.7$, $\Delta G_s = +2.12$ kcal/mol, and $\Delta G_{Dangle} = +1.54$ kcal/mol.

Optical Melting Measurements on Helices with D-U Base pairs

To test the prediction accuracy of RECCES, we experimentally measured the folding free energies of RNA helices containing D-U, a non-natural base pair for which the NN parameters have not been previously determined. We performed optical melting experiments on six helices containing D-U using a Shimadzu Spectrophotometer UV-1800, as well as two helices with canonical base pairs, which were confirmed to give thermodynamic parameters that reproduced literature values. The RNA constructs were ordered from Dharmacon with HPLC purification. For each construct, we obtained 12 melting curves at three different concentrations (15, 25 and 35 μM), measured at 260 nm. The buffer system was 1.0 M NaCl, 20 mM sodium cacodylate at pH 7.0, and 0.5 mM Na_2EDTA . For each sample, RNA concentrations were determined *in situ* using the high-temperature absorbances and extinction coefficients of RNA single-strands(18, 19). The extinction coefficient parameter for 2,6-diaminopurine is unknown. The free energies computed herein are insensitive to small changes of the extinction coefficients (20); we assumed that 2,6-diaminopurine has the same parameters as adenine. The melting curves were fitted with a two-state model with linear baselines to obtain the corresponding folding free energies for all constructs (21-23). To obtain the NN parameters from the helix folding free energies, recall that the helix folding free energies can be written as linear combinations of NN parameters (e.g., see eqs. (1), (3) and (5)). The NN parameters can be solved through least-squares regression, with the exception of $\Delta G_{\text{Terminal-DU}}$, the terminal contribution of D-U relative to G-C. Determining $\Delta G_{\text{Terminal-DU}}$ requires measurements on RNA sequences with D at the 5'-end of strands and with D on the 3'-end of the strands (see, e.g., eq. (4)), but the latter are not commercially available. For results in the main text, we have assumed that this parameter is zero, as it should have the same geometry and same number of hydrogen-bonds (3) as G-C. Assuming different terminal D-U contributions from -0.2 to $+0.3$ kcal/mol

(compiled in Supporting Information) gives similar RMSD accuracies of our predictions compared to experimental values.

Supporting Results

Accuracies of terminal base pair penalty calculations

In addition to predicting NN parameters for two-base-pair segments, we were able to calculate terminal contributions with RECCES-Rosetta (see main text, eq. (4)). The terminal contribution is the free energy contribution for having a non-G-C base pair instead of a G-C pair at the end of a helix. For example, because an A-U base pair only has two hydrogen bonds, one less than that of G-C pair, putting it at the end of the helix incurs a penalty in folding free energy. Our method predicted the terminal contribution for A-U pair to be 0.75 ± 0.15 kcal/mol, somewhat higher than the experimental value (0.45 kcal/mol)(21). For G-U pairs, our prediction suggested it to be similar but slightly less stable than A-U pairs at helix terminal (0.87 ± 0.3 kcal/mol), while recent experiments have shown that it is equally stable as G-C pair at terminal (24) (0 terminal contribution). The higher experimental stability of G-U pair at the terminal may be due to the variety of non-A-form alternative conformations that it is known to the sample (24); our current calculations only sampled near-A-form conformations for helices. For the iG-iC pair, our method predicted a terminal contribution of -0.11 ± 0.16 kcal/mol, within error of the experimental value (-0.19 ± 0.07 kcal/mol) (20).

Supporting Tables

Table S1. Systematic errors due to different system setups, sampling schemes, and parameters used in WHAM analysis.

| | Original | No initial BP ¹ | Initial BP = A/U ¹ | Initial BP = GG/CC ¹ | Reduced A-form range ² | Alternative Sugar ³ | Modified bin-size and cutoff ⁴ |
|-------------------------------|----------|----------------------------|-------------------------------|---------------------------------|-----------------------------------|--------------------------------|---|
| 5'AC3' 3'UG5' | -1.79 | -1.91 | -1.62 | -1.49 | -1.70 | -1.72 | -1.79 |
| 5'AG3' 3'UC5' | -2.14 | -2.17 | -1.83 | -1.92 | -2.08 | -2.31 | -2.14 |
| 5'UC3' 3'AG5' | -2.10 | -2.32 | -1.94 | -1.82 | -2.00 | -1.92 | -2.10 |
| 5'UG3' 3'AC5' | -1.83 | -1.46 | -1.52 | -1.58 | -1.89 | -1.91 | -1.83 |
| Error (kcal/mol) ⁵ | | 0.23 | 0.26 | 0.27 | 0.08 | 0.14 | 5.86×10 ⁻⁵ |

¹ The initial base pairs for the simulation (see methods).

² Use ±40° as the angle constraint for A-form helix instead of ±60°.

³ Use alternative conformation for the sugar rings during sampling (“-rna::corrected_geo false” in Rosetta command line).

⁴ Use different bin-size and high-score cutoff (0.05 and 1600) instead. The original value is 0.1 and 800.

⁵ Error relative to the original.

Table S2. Statistical error on sampling computed based on bootstrapping (resampling of model ensembles with replacement).

| Sequence | Simulated molecule | Mean free energy (kcal/mol) | Standard deviation (kcal/mol) |
|--------------------|-----------------------|-----------------------------|-------------------------------|
| 5'GGG3' 3'CCC5' | Complex | -5.54 | 0.0023 |
| 5'GU3' 3'CA5' | Complex | -5.12 | 0.0010 |
| 5'ACC3' | Single strand | 10.86 | 0.0017 |
| 5'GU3' | Single strand | 5.60 | 0.0014 |
| 5'GC3' 3'CGA5' | Complex, dangling end | -5.91 | 0.0010 |

Table S3. Systematic error from RECCES reweighting procedure.

NN predictions from five weight sets were computed by applying a fast reweighting procedure to models simulated with an arbitrarily chosen starting weight set, and compared to predictions from explicit simulation of the conformational ensemble with five arbitrarily chosen new weight sets.

| Score term | Starting | Weight 1 | Weight 2 | Weight 3 | Weight 4 | Weight 5 |
|-------------------------------------|----------|----------|----------|----------|----------|----------|
| <i>fa_atr</i> | 0.37 | 0.58 | 0.28 | 0.10 | 0.33 | 0.93 |
| <i>fa_rep</i> | 0.2 | 0.1 | 1.25 | 0.72 | 0.32 | 1.15 |
| <i>fa_intra_rep</i> | 0.0035 | 0.0022 | 0.0046 | 0.0018 | 0.01 | 0 |
| <i>fa_stack</i> | 0.00001 | 0.0203 | 6.00 | 0.40 | 0.25 | 0 |
| <i>rna_torsion</i> | 9.5 | 4.90 | 5.79 | 5.65 | 3.29 | 6.78 |
| <i>hbond_sc</i> | 3.8 | 2.61 | 0.61 | 5.39 | 3.30 | 5.85 |
| <i>lk_nonpolar</i> | 0.51 | 0 | 0.61 | 1.83 | 0.038 | 0.92 |
| <i>geom_sol_fast</i> | 0.40 | 0 | 1.68 | 1.40 | 0 | 2.02 |
| <i>stack_elec</i> | 1.7 | 3.26 | 2.29 | 2.41 | 1.92 | 1.14 |
| <i>fa_elec_rna_phos_phos</i> | 1.6 | 6.77 | 1.34 | 0 | 4.47 | 0.72 |
| Reweighting error (kcal/mol) | n.a. | 0.48 | 0.26 | 0.25 | 0.22 | 0.23 |

Table S4. Mean and standard deviation of the weights of each score term from different optimization runs.

| Score term | Description | Mean | Standard deviation |
|------------------------------|---|-------|--------------------|
| <i>fa_atr</i> | Lennard-Jones attraction | 0.52 | 0.23 (45%) |
| <i>fa_rep</i> | Lennard-Jones repulsion (inter-residue) | 0.27 | 0.51 (186%) |
| <i>fa_intra_rep</i> | Lennard-Jones repulsion (intra-residue) | 0.067 | 0.038 (57%) |
| <i>fa_stack</i> | Extra Lennard-Jones attraction for stacking atoms | 0.044 | 0.080 (181%) |
| <i>rna_torsion</i> | RNA torsional potential | 9.3 | 4.7 (50%) |
| <i>hbond_sc</i> | Hydrogen bond | 3.8 | 1.6 (42%) |
| <i>lk_nonpolar</i> | Lazaridis-Karplus solvation (Ref. (25)) | 1.1 | 0.94 (85%) |
| <i>geom_sol_fast</i> | Geometric solvation (Ref. (10)) | 0.83 | 1.2 (139%) |
| <i>stack_elec</i> | Electrostatics for stacking atoms | 1.4 | 1.0 (74%) |
| <i>fa_elec_rna_phos_phos</i> | Electrostatics for backbone phosphates | 9.6 | 7.6 (80%) |

Table S5. Example component weights in RECCES collection of energy functions.

| Score term | Structure prediction ¹ | Best ² | Weight 1 | Weight 2 | Weight 3 | Weight 4 |
|--|-----------------------------------|-------------------|----------|----------|----------|----------|
| <i>fa_atr</i> | 0.23 | 0.73 | 0.72 | 0.94 | 0.50 | 0.39 |
| <i>fa_rep</i> | 0.12 | 0.1 | 0.17 | 0.39 | 0.1 | 0.19 |
| <i>fa_intra_rep</i> | 0.0029 | 0.0071 | 0.01 | 0.0029 | 0.0042 | 0.0100 |
| <i>fa_stack</i> | 0 | 0 | 0 | 0 | 0 | 0.051 |
| <i>rna_torsion</i> | 0.1 | 4.26 | 4.30 | 3.46 | 8.45 | 18.64 |
| <i>hbond_sc</i> | 3.4 | 2.46 | 2.48 | 2.36 | 3.11 | 5.99 |
| <i>lk_nonpolar</i> | 0.32 | 0.25 | 1.66 | 0.82 | 0.55 | 3.63 |
| <i>geom_sol_fast</i> | 0.62 | 0 | 0 | 0 | 0.38 | 2.51 |
| <i>stack_elec</i> | 0 | 1.54 | 0.77 | 1.15 | 0.75 | 1.49 |
| <i>fa_elec_rna_phos_phos</i> | 1.05 | 4.54 | 0 | 0.73 | 13.6 | 20 |
| Training error (kcal/mol) | 1.23 | 0.35 | 0.38 | 0.41 | 0.36 | 0.38 |
| Test error, GU (kcal/mol) | 0.93 | 0.34 | 0.43 | 0.38 | 0.34 | 0.51 |
| Test error, iGiC (kcal/mol) | 1.79 | 1.06 | 0.96 | 1.10 | 1.00 | 1.03 |
| Test error, iGiC with outlier exclusion (kcal/mol) | 1.26 | 0.52 | 0.54 | 0.69 | 0.62 | 0.66 |

¹ The default score weights used in Rosetta structure prediction. In addition to the listed score terms, it also includes the terms *ch_bond* (0.42), *hbond_sr_bb_sc* (0.62), *hbond_lr_bb_sc* (3.4) (weights in parentheses). Here *ch_bond* is the hydrogen bond interaction between C-H to polar atom (O and N). The other two terms are hydrogen bonds between side-chain and backbone atoms. The terms do not contribute in a sequence-dependent manner to helix NN parameters and were omitted from the simulations herein. The training and test errors presented in the table are based on calculations optimizing a global scaling factor converting Rosetta energy units to $k_B T$ (final value of $1 k_B T / 0.218$) over training data (canonical stacked pairs and dangling ends).

² The best weight with the lowest training error in the reweighting.

Table S6. Highly correlated score terms in the collection of RECCES score functions.

Each score function uses a different, locally optimized weight set fitted to canonical Watson-Crick stacked pairs and dangling ends for natural bases. Terms with absolute correlation > 0.1 are listed.

| Score term 1 | Score term 2 | Correlation coefficient |
|----------------------|------------------------------|-------------------------|
| <i>fa_atr</i> | <i>fa_rep</i> | 0.31 |
| <i>fa_atr</i> | <i>fa_stack</i> | -0.62 |
| <i>fa_atr</i> | <i>rna_torsion</i> | -0.35 |
| <i>fa_atr</i> | <i>stack_elec</i> | -0.11 |
| <i>fa_atr</i> | <i>fa_elec_rna_phos_phos</i> | -0.17 |
| <i>fa_rep</i> | <i>rna_torsion</i> | -0.12 |
| <i>fa_rep</i> | <i>lk_nonpolar</i> | 0.11 |
| <i>fa_rep</i> | <i>stack_elec</i> | 0.14 |
| <i>fa_rep</i> | <i>fa_elec_rna_phos_phos</i> | -0.18 |
| <i>fa_stack</i> | <i>rna_torsion</i> | -0.12 |
| <i>fa_stack</i> | <i>lk_nonpolar</i> | -0.19 |
| <i>fa_stack</i> | <i>stack_elec</i> | 0.35 |
| <i>fa_stack</i> | <i>fa_elec_rna_phos_phos</i> | -0.10 |
| <i>rna_torsion</i> | <i>hbond_sc</i> | 0.73 |
| <i>rna_torsion</i> | <i>lk_nonpolar</i> | 0.47 |
| <i>rna_torsion</i> | <i>geom_sol_fast</i> | 0.72 |
| <i>rna_torsion</i> | <i>fa_elec_rna_phos_phos</i> | 0.61 |
| <i>hbond_sc</i> | <i>lk_nonpolar</i> | 0.29 |
| <i>hbond_sc</i> | <i>geom_sol_fast</i> | 0.99 |
| <i>hbond_sc</i> | <i>stack_elec</i> | 0.32 |
| <i>hbond_sc</i> | <i>fa_elec_rna_phos_phos</i> | 0.63 |
| <i>lk_nonpolar</i> | <i>geom_sol_fast</i> | 0.32 |
| <i>lk_nonpolar</i> | <i>fa_elec_rna_phos_phos</i> | 0.26 |
| <i>geom_sol_fast</i> | <i>stack_elec</i> | 0.29 |
| <i>geom_sol_fast</i> | <i>fa_elec_rna_phos_phos</i> | 0.64 |

Table S7. Experimental measurements and free energy predictions of helices containing 2,6-diaminopurine (D)-uracil base pairs.

| Sequence | Experiment | RECCES-Rosetta | Hydrogen-bond counting | Rosetta score |
|-----------------|--------------|----------------|------------------------|---------------|
| GDGCUC | -9.56 ± 0.37 | -11.51 ± 0.38 | -10.04 ± 0.24 | -10.34 ± 0.24 |
| CUGCDG | -9.06 ± 0.38 | -12.45 ± 0.42 | -11.11 ± 0.23 | -10.61 ± 0.23 |
| CDCGUG | -8.34 ± 0.11 | -10.58 ± 0.44 | -10.04 ± 0.24 | -10.04 ± 0.24 |
| DGCGCU | -9.16 ± 0.63 | -11.14 ± 0.42 | -10.84 ± 0.38 | -10.46 ± 0.38 |
| GDCGUC | -9.38 ± 0.35 | -9.64 ± 0.4 | -11.11 ± 0.23 | -10.91 ± 0.23 |
| DCCGGU | -9.67 ± 0.22 | -10 ± 0.43 | -10.5 ± 0.37 | -10.62 ± 0.37 |
| RMSE (kcal/mol) | | 2.02 (21.9%) | 1.52 (16.5%) | 1.34 (14.6%) |

Predicted values are computed by linear combinations of the predicted NN parameters; errors for the predictions were estimated by error propagation.

Table S8. Nearest neighbor parameters for 2,6-Diaminopurine/U containing stacked pairs assuming different terminal D-U contributions (in kcal/mol).

| Terminal D-U | 5'DG/3'UC | 5'GD/3'CU | 5'DC/3'UG | 5'CD/3'GU | RMSE (vs. RECCES) |
|----------------|-----------|-----------|-----------|-----------|-------------------|
| -0.2 | -2.08 | -3.30 | -2.42 | -2.92 | 0.68 |
| -0.1 | -2.18 | -3.20 | -2.52 | -2.82 | 0.65 |
| 0 | -2.28 | -3.10 | -2.62 | -2.72 | 0.63 |
| 0.1 | -2.38 | -3.00 | -2.72 | -2.62 | 0.63 |
| 0.2 | -2.48 | -2.90 | -2.82 | -2.52 | 0.65 |
| 0.3 | -2.58 | -2.80 | -2.92 | -2.42 | 0.67 |
| RECCES-Rosetta | -3.2 | -3.1 | -2.8 | -3.57 | n.a. |

Supporting Figures

Figure S1. Correlation plots for the component score terms.

Each red cross represents a minimized score weight set from the RECCES-Rosetta energy function collection. The blue lines are linear regressions of the red crosses.

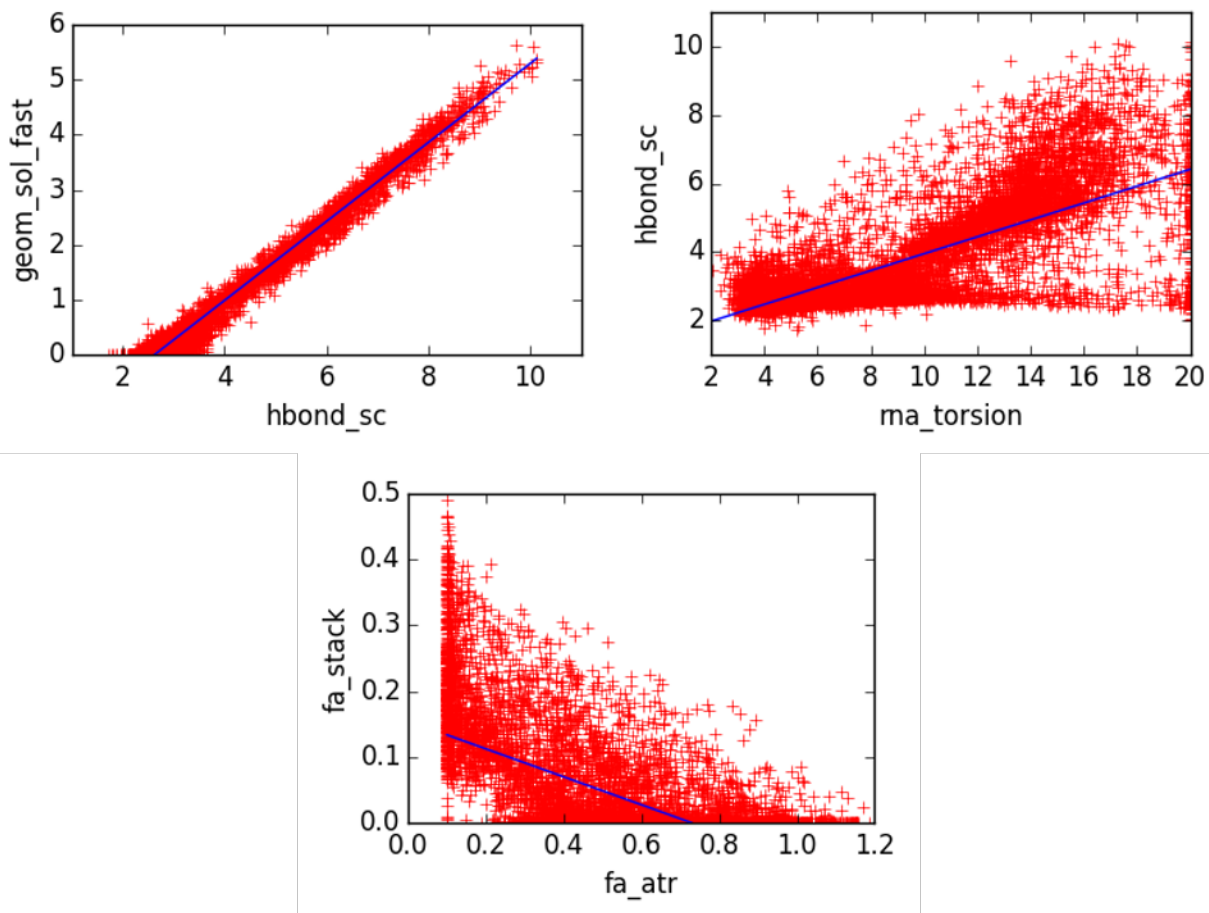
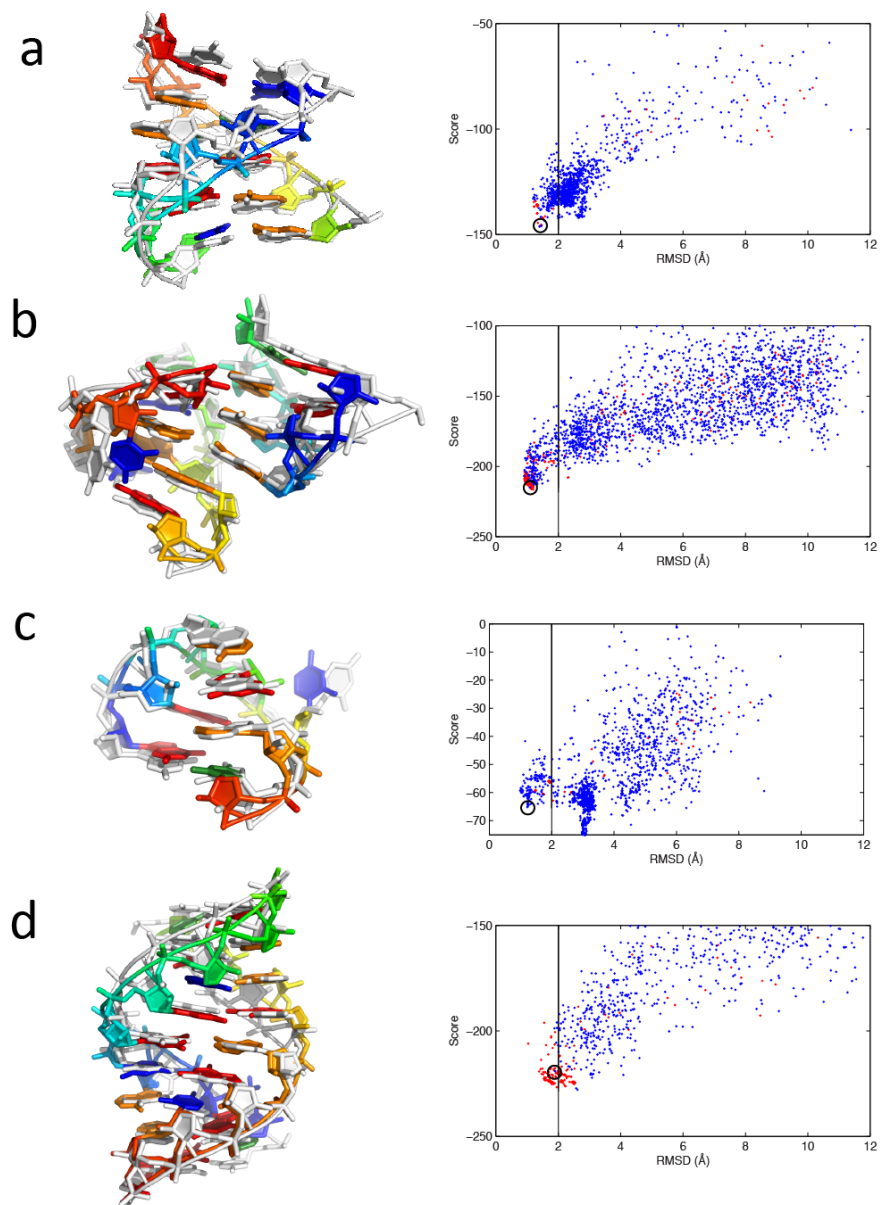


Figure S2. RECCES weights tested in structure prediction.

De novo 3D structure prediction (5000 fragment assembly steps, 2000 models each) repeated for all noncanonical motifs in the FARFAR benchmark (Das, Karanicolas, Baker, Nature Methods, 2010) confirms that the ‘best’ RECCES weight set (Table S5) recovers 12 of the 16 motifs previously recovered at better than 2 Å resolution. Examples include (a) domain IV from signal recognition particle RNA (1LNT) and (b) the tetraloop/receptor from the P4-P6 domain (2R8S). For the other 4 motifs, models better than 2 Å accuracy are still sampled, are stable to refinement in the RECCES energy function, and give energies similar to the lowest energy models. Examples include (a) GAGUA pentaloop from a SARS conserved domain (1XJR) and (b) the bacterial 5S ribosomal RNA loop E motif. Each panel gives single-model RECCES energies for *de novo* models (blue) and refined native models (red) vs. RMSD to crystallographic structure (right); and one 3D model (colored) overlaid on crystallographic structure (white) (on left). These results are promising; the RECCES scores for single motif conformations should be less accurate than free energies estimated for each conformation, whose computation may be possible with extensions of RECCES and would also allow inclusion of structure prediction during weight training.



References for Supporting Information Appendix

1. Serra MJ & Turner DH (1995) Predicting thermodynamic properties of RNA. *Meth. Enzymol.* 259:242-261.
2. Das R, Karanicolas J, & Baker D (2010) Atomic accuracy in predicting and designing noncanonical RNA structure. *Nat Methods* 7(4):291-294.
3. Marinari E & Parisi G (1992) Simulated Tempering: A New Monte Carlo Scheme. *EPL* 19(6):451.
4. Huang X, Bowman GR, & Pande VS (2008) Convergence of folding free energy landscapes via application of enhanced sampling methods in a distributed computing environment. *The Journal of Chemical Physics* 128(20):205106.
5. Kumar S, Rosenberg J, Bouzida D, Swendsen R, & Kollman P (1992) The weighted histogram analysis method for free-energy calculations on biomolecules. I. The method. *J. Comput. Chem.* 13(8):1011-1021.
6. Chodera JD, Swope WC, Pitera JW, Seok C, & Dill KA (2006) Use of the Weighted Histogram Analysis Method for the Analysis of Simulated and Parallel Tempering Simulations. *J. Chem. Theory Comput.* 3(1):26-41.
7. Guerois R, Nielsen JE, & Serrano L (2002) Predicting Changes in the Stability of Proteins and Protein Complexes: A Study of More Than 1000 Mutations. *Journal of Molecular Biology* 320(2):369-387.
8. Leaver-Fay A, *et al.* (2013) Chapter Six - Scientific Benchmarks for Guiding Macromolecular Energy Function Improvement. *Meth. Enzymol.*, Methods in Protein Design, ed Keating AE (Academic Press), Vol 523, pp 109-143.
9. Park S & Pande VS (2007) Choosing weights for simulated tempering. *Phys. Rev. E* 76(1).
10. Das R, Karanicolas J, & Baker D (2010) Atomic accuracy in predicting and designing noncanonical RNA structure. *Nat Meth* 7(4):291-294.
11. Morozov AV, Kortemme T, Tsemekhman K, & Baker D (2004) Close agreement between the orientation dependence of hydrogen bonds observed in protein structures and quantum mechanical calculations. *PNAS* 101(18):6946-6951.
12. Sripakdeevong P, Kladwang W, & Das R (2011) An enumerative stepwise ansatz enables atomic-accuracy RNA loop modeling. *Proceedings of the National Academy of Sciences* 108(51):20573-20578.
13. Chou F-C, Sripakdeevong P, Dibrov SM, Hermann T, & Das R (2013) Correcting pervasive errors in RNA crystallography through enumerative structure prediction. *Nat Meth* 10(1):74-76.
14. Sripakdeevong P, *et al.* (2014) Structure determination of noncanonical RNA motifs guided by 1H NMR chemical shifts. *Nat Meth* 11(4):413-416.
15. O'Meara MJ, *et al.* (2015) A Combined Covalent-Electrostatic Model of Hydrogen Bonding Improves Structure Prediction with Rosetta. *J Chem Theory Comput* 11(2):609-622.
16. Yildirim I & Turner DH (2005) RNA Challenges for Computational Chemists†. *Biochemistry* 44(40):13225-13234.
17. Yesselman JD, Price DJ, Knight JL, & Brooks CL (2012) MATCH: An atom-typing toolset for molecular mechanics force fields. *Journal of Computational Chemistry* 33(2):189-202.
18. Cantor CR & Tinoco Jr I (1965) Absorption and Optical Rotatory Dispersion of Seven Trinucleoside Diphosphates. *Journal of Molecular Biology* 13(1):65-77.
19. Murugaiah V (2011) Determination of Extinction Coefficient. *Handbook of Analysis of Oligonucleotides and Related Products*, ed Srivatsa G (CRC Press), pp 351-358.

20. Chen X, Kierzek R, & Turner DH (2001) Stability and Structure of RNA Duplexes Containing Isoguanosine and Isocytidine. *Journal of the American Chemical Society* 123(7):1267-1274.
21. Xia T, *et al.* (1998) Thermodynamic Parameters for an Expanded Nearest-Neighbor Model for Formation of RNA Duplexes with Watson–Crick Base Pairs. *Biochemistry* 37(42):14719-14735.
22. Mathews DH, Sabina J, Zuker M, & Turner DH (1999) Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure. *Journal of Molecular Biology* 288(5):911-940.
23. Andronescu M, Condon A, Turner DH, & Mathews DH (2014) The Determination of RNA Folding Nearest Neighbor Parameters. *RNA Sequence, Structure, and Function: Computational and Bioinformatic Methods*, Methods in Molecular Biology, eds Gorodkin J & Ruzzo WL (Humana Press), pp 45-70.
24. Chen JL, *et al.* (2012) Testing the Nearest Neighbor Model for Canonical RNA Base Pairs: Revision of GU Parameters. *Biochemistry* 51(16):3508-3522.
25. Lazaridis T & Karplus M (1999) Effective energy function for proteins in solution. *Proteins: Structure, Function, and Bioinformatics* 35(2):133-152.