# Three-dimensional RNA structure refinement by hydroxyl radical probing

Feng Ding[1,2], Christopher A Lavender[3], Kevin M Weeks[3] & Nikolay V Dokholyan[1,2]

**Molecular modeling guided by experimentally derived structural information is an attractive approach for three-dimensional structure determination of complex RNAs that are not amenable to study by high-resolution methods. Hydroxyl radical probing (HRP), which is performed routinely in many laboratories, provides a measure of solvent accessibility at individual nucleotides. HRP measurements have, to date, only been used to evaluate RNA models qualitatively. Here we report the development of a quantitative structure refinement approach using HRP measurements to drive discrete molecular dynamics simulations for RNAs ranging in size from 80 to 230 nucleotides. We first used HRP reactivities to identify RNAs that form extensive helical packing interactions. For these RNAs, we achieved highly significant structure predictions given the inputs of RNA sequence and base pairing. This HRP-directed tertiary structure refinement approach generates robust structural hypotheses that are useful for guiding explorations of structure-function inter-relationships in RNA.**

RNA molecules have central roles in gene expression, splicing and translation[1]. Knowledge of their underlying three-dimensional structure is a fundamental prerequisite to a complete understanding of most RNA functions. High-resolution methods such as X-ray crystallography and nuclear magnetic resonance spectroscopy offer unparalleled atomic-level insight into RNA structure. However, many RNAs are not amenable to structural characterization by these methods because of their conformational flexibility or large size. Recent advances[2–5] in molecular modeling have yielded accurate structure predictions of small RNAs, but because of the vast RNA conformational space and because the available force fields do not accurately describe atomic interactions, structure prediction for large RNA molecules with complex topologies is beyond the reach of the current *ab initio* approaches. Incorporation of experimentally derived structural information with computational modeling can markedly reduce the allowed conformational space and thereby facilitate the prediction of native RNA ensembles[6–11].

The pattern of base pairing in an RNA, also known as its secondary structure, can often be established with high accuracy by sequence covariation analysis[12,13] or experimentally constrained secondary structure prediction, especially when using information obtained from selective 2′-hydroxyl acylation analyzed by primer extension experiments[14,15]. Accurate knowledge of the secondary structure of an RNA greatly restrains the tertiary folds that are possible[16,17], but the size of the conformational space is still large[16]. Through-space distance constraints derived from biochemical experiments or bioinformatics analyses can provide information that is crucial for refining the folds of an RNA molecule. A small number of long-range, through-space distance constraints are often sufficient to limit the conformational space enough to allow accurate RNA structure prediction[10,12]. Experimental methods that are used to probe through-space distances, including site-directed hydroxyl radical footprinting, crosslinking and fluorescence resonance energy transfer, can give high-quality distance information. However, these approaches often require the synthesis of specialized RNA constructs, careful controls for unintended structural perturbations and complex approaches for data interpretation[15]. In contrast, HRP, which reports the approximate backbone solvent accessibility[18–20] (**Fig. 1a**), is relatively straightforward to implement. HRP measurements have previously been used to evaluate or filter RNA structural ensembles[9,18,21,22] but not to drive three-dimensional RNA structure determination in a quantitative and systematic way. Here we describe a framework for biasing discrete molecular dynamics (DMD)[23] simulations of RNA to generate structural ensembles that are consistent with experimental HRP measurements.
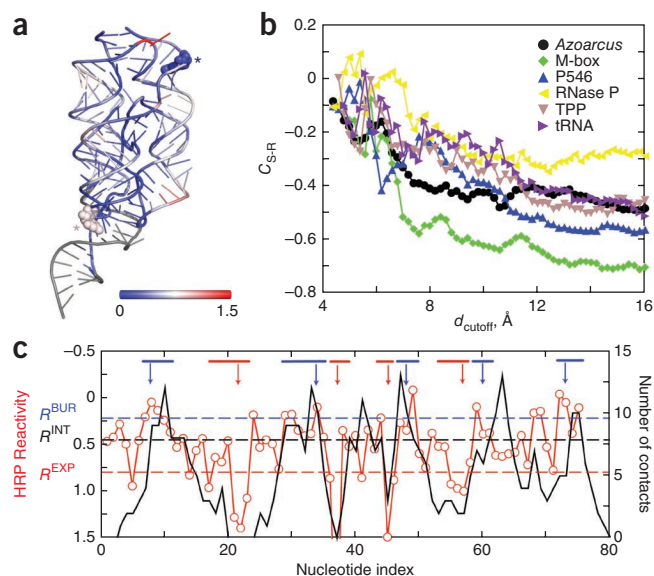
## RESULTS

### A coarse-grained RNA model for driving RNA structure refinement

We used a coarse-grained approach to model RNA molecules in which each RNA nucleotide is represented by three pseudoatoms corresponding to the base, sugar and phosphate groups. Three beads are sufficient to correctly recapitulate the major features of the RNA structure, including excluded volume, base pairing and stacking and loop entropy, and this model is sufficiently simple to allow for efficient computational sampling[3]. This three-bead modeling approach has been used successfully to fold small RNAs with simple topologies from their sequence alone[3] and to refine

[1]Department of Biochemistry and Biophysics, University of North Carolina, Chapel Hill, North Carolina, USA. [2]Center for Computational and Systems Biology, University of North Carolina, Chapel Hill, North Carolina, USA. [3]Department of Chemistry, University of North Carolina, Chapel Hill, North Carolina, USA. Correspondence should be addressed to N.V.D. (dokh@unc.edu) or K.M.W. (weeks@unc.edu).

**Figure 1** | The relationship between RNA structure and HRP reactivity. (**a**) The structure of the M-box riboswitch is illustrated. Nucleotides are colored according to HRP reactivity (blue to red); nucleotides without HRP data are shown in gray. A solvent-exposed nucleotide with low HRP reactivity (blue) and a buried nucleotide with high HRP reactivity (red) are emphasized with all-atom representations (asterisks). (**b**) The structure-reactivity correlation coefficient, $C_{\text{S-R}}$, as a function of the $d_{\text{cutoff}}$ values for the six training RNAs using HRP data smoothed over a 3-nt window (Online Methods). (**c**) Comparison of experimentally measured HRP reactivities (red) with the number of through-space contacts (black) for the TPP riboswitch RNA using a $d_{\text{cutoff}}$ of 14.0 Å. Buried and exposed nucleotide segments are denoted with blue and red lines, respectively (top); the downward colored arrows indicate the representative nucleotides that are characteristic of each nucleotide segment. The dashed horizontal lines represent the exposed ($R^{\text{EXP}}$), buried ($R^{\text{BUR}}$) and intermediate ($R^{\text{INT}}$) threshold values.



larger RNA structures using distance constraints[8,10]. This delineation into the base, sugar and phosphate groups is also compatible with HRP chemistry, where the hydroxyl radical reacts primarily at the ribose sugar[18].

We first optimized HRP-directed refinement with a training set of six structurally diverse RNAs ranging from 75 to 214 nucleotides in length (**Table 1**). We evaluated the prediction accuracy by comparison with available high-resolution structures. After optimization with the training set, we applied HRP-directed refinement to an independent set of four RNAs (from 152 to 412 nucleotides in length; **Table 1**). We did not use structures in the test set to optimize the method, and, therefore, we expect the accuracy of the resulting models to be indicative of the predictive capability of the HRP-directed structure refinement method.

## Correlation of HRP reactivity with number of neighbor atoms

We used a widely applied approach for the HRP experiment[24]. Our data are consistent with protection patterns described in previous HRP experiments (Online Methods). To incorporate experimentally measured HRP reactivities into the DMD simulations, we needed to define a structural parameter that was reflective of the information obtained in an HRP measurement and that could also be readily implemented as a constraint to drive the folding simulations. Hydroxyl radical reactivity is correlated with backbone solvent accessibility[19,25]; however, it is not straightforward to incorporate solvent accessibility as a constraint in a molecular dynamics simulation. We found that solvent accessibility is inversely proportional to the number of through-space neighboring nucleotides, which we term contacts (**Fig. 1**). In the example of the M-box riboswitch, with the exception of some outliers (**Fig. 1a**), nucleotides with low HRP reactivities were generally buried and had many through-space contacts, whereas nucleotides with high reactivities had fewer through-space contacts and were more exposed (**Fig. 1a**). The number of through-space contacts can be readily incorporated as a constraint in DMD and other simulation methods, and we used it here to bias our simulations (Online Methods).

We defined through-space contacts based on the sugar pseudo-atoms in our three-bead model of RNA[3]. We computed the number of contacts as the number of sugar beads within a cutoff distance, $d_{\text{cutoff}}$, of a given nucleotide sugar bead. We excluded immediate neighbors in the linear sequence and base-pairing partners in helical elements because these neighbors reflect

primary and secondary structure rather than higher-order tertiary interactions. To find the optimal $d_{\text{cutoff}}$ values, we calculated the structure-reactivity correlation, $C_{\text{S-R}}$, as the Pearson correlation coefficient between the number of contacts and the corresponding HRP reactivity for each nucleotide (**Fig. 1b**). We used the data obtained from the six RNAs from the training set to determine the optimal $d_{\text{cutoff}}$ values. In these calculations, $C_{\text{S-R}}$ was negative because a lower HRP reactivity corresponds to a more buried nucleotide with a larger number of through-space neighbors. The absolute magnitude of $C_{\text{S-R}}$ was largest when the $d_{\text{cutoff}}$ range was 13–15 Å (**Fig. 1b**). With an intermediate $d_{\text{cutoff}}$ value of 14 Å, the correlation coefficients for the six training RNAs ranged from −0.5 to −0.7, with the exception of RNase P, for which $C_{\text{S-R}}$ was smaller (approximately −0.30).

## Interaction potentials assigned according to HRP reactivities

To incorporate the HRP data into the DMD simulations, we assigned two bias interaction potentials (Online Methods). The first included pairwise attractive potentials for all nucleotides to encourage the collapse of the RNA and general nucleotide packing. The second was an overburial repulsion potential incurred when a given nucleotide exceeded the assigned threshold number of contacts ($N_{\text{max}}$) derived from its experimental HRP reactivity (**Fig. 2a**). To assign $N_{\text{max}}$ values to the nucleotides, we defined high and low HRP cutoff values corresponding to the highest and lowest mean HRP reactivities, respectively. Based on an analysis of single-chain RNAs in the Protein Data Bank[26] and on exploratory simulations with the six training set RNAs, we assigned the largest and smallest $N_{\text{max}}$ values as 11 and 0.5, respectively (**Fig. 2b**). We assigned the $N_{\text{max}}$ values of 11 and 0.5 to the nucleotides with HRP values above and below the HRP threshold, respectively. For nucleotides with intermediate HRP values, we assigned $N_{\text{max}}$ values using linear interpolation (Online Methods).

HRP experiments are intrinsically noisy (**Fig. 1c** and **Supplementary Fig. 1**), making assignment of interaction potentials challenging, especially in regions with moderate HRP reactivities. To reduce the effects of noise on structure prediction, we incorporated stronger-biasing interactions for RNA nucleotides that we could designate as exposed or buried with high confidence.

**Table 1** | Summary of HRP-directed RNA fold refinement for the studied RNAs

| RNA | Length (nt) | $f_{0.25}$ | $C_{S-R}$ | Number of clusters | Large clusters n (out of 100) | RMSD (Å) | P value |
|---|---|---|---|---|---|---|---|
| *Azoarcus* group I intron | 214 | **0.43** | −0.45 | 1 | 100 | 16.8 ± 2.1 | $<10^{-6}$ |
| M-box riboswitch | 161 | **0.35** | −0.68 | 1 | 100 | 11.6 ± 2.3 | $<10^{-6}$ |
| P546 domain | 158 | **0.37** | −0.57 | 3 | 66 | 19.8 ± 1.4 | $3.0 \times 10^{-3}$ |
| | | | | | 32 | 15.1 ± 1.9 | $<10^{-6}$ |
| RNase P specificity domain | 152 | 0.25 | −0.30 | 3 | 93 | 24.9 ± 1.2 | 0.67 |
| | | | | | 4 | 24.1 ± 3.2 | 0.50 |
| | | | | | 3 | 22.7 ± 0.5 | 0.22 |
| TPP riboswitch | 80 | 0.21 | −0.50 | 4 | 44 | 12.6 ± 1.3 | 0.10 |
| | | | | | 42 | 14.6 ± 1.9 | 0.44 |
| | | | | | 12 | 9.9 ± 1.8 | $2.9 \times 10^{-3}$ |
| tRNA$^{Asp}$ | 75 | 0.25 | −0.45 | 8 | 29 | 14.1 ± 1.3 | 0.50 |
| | | | | | 23 | 17.5 ± 1.5 | 0.97 |
| | | | | | 18 | 16.4 ± 0.8 | 0.90 |
| | | | | | 8 | 18.9 ± 0.8 | 0.99 |
| | | | | | 6 | 12.3 ± 1.5 | 0.17 |
| | | | | | 6 | 15.7 ± 0.4 | 0.82 |
| | | | | | 5 | 18.9 ± 0.8 | 0.99 |
| *O. iheyensis* group II intron | 412 | 0.21 | −0.30 | – | – | – | – |
| RNase P catalytic domain | 231 | **0.28** | −0.50 | 4 | 46 | 19.2 ± 1.6 | $<10^{-6}$ |
| | | | | | 41 | 21.6 ± 2.1 | $<10^{-6}$ |
| | | | | | 8 | 25.0 ± 0.7 | $1.3 \times 10^{-4}$ |
| | | | | | 5 | 24.4 ± 2.0 | $3.5 \times 10^{-5}$ |
| Lysine riboswitch | 174 | **0.36** | −0.57 | 3 | 57 | 12.0 ± 1.6 | $<1.0^{-6}$ |
| | | | | | 42 | 18.1 ± 1.0 | $1.6 \times 10^{-6}$ |
| glmS ribozyme | 152 | **0.35** | −0.55 | 2 | 74 | 16.6 ± 1.7 | $1.5 \times 10^{-5}$ |
| | | | | | 26 | 8.5 ± 1.3 | $<10^{-6}$ |

The first six RNAs listed comprise the training set used for the algorithm optimization and applicability determination: the yeast tRNA$^{Asp}$ (ref. 27), TPP riboswitch[28], specificity domain of RNase P[29], P546 domain of the *T. thermophila* group I intron[19], M-box riboswitch[24] and *Azoarcus* group I intron[30] RNAs. The last four RNAs listed were used for testing the performance of our HRP-directed fold refinement method: the glmS ribozyme[31], lysine riboswitch[32], catalytic domain of RNase P[33] and *O. iheyensis* group II intron[34] RNAs. The fraction of highly protected nucleotides, $f_{0.25}$, was computed using only the experimental HRP data; $f_{0.25}$ values above 0.25 are shown in bold. The structure-reactivity correlation, $C_{S-R}$, was calculated with reference to the accepted experimental structure. One hundred selected structures were clustered by pairwise RMSD (Online Methods). Small clusters (with one or two structures) were excluded. For each cluster, the P values were calculated based on the average RMSD with respect to the accepted experimental structure[16]. HRP-directed structure refinement was not performed for the *O. iheyensis* group II intron RNA in the test set because of its low $f_{0.25}$ value.

We identified RNA segments (≥3 nt) with high or low HRP reactivities and selected the central nucleotide in each segment as the representative exposed or buried nucleotide, respectively (Online Methods and **Fig. 1c**, red and blue bars). These central representative nucleotides had a high probability of being buried or exposed in the native structure because the effect of the noise associated with the HRP measurements is less pronounced when measured over several consecutive highly buried or solvent-exposed nucleotides. We included a strong pairwise attraction between nucleotides identified as highly buried and the rest of the RNA molecule, whereas we assigned a strong overburial repulsion for the nucleotides identified as either highly buried or highly exposed (Online Methods).

### HRP-directed RNA structure refinement for the tested RNAs

We used DMD simulations in three steps to obtain structural ensembles that were consistent with the experimental HRP data (**Fig. 2c** and Online Methods). First, we performed serial DMD simulations with inputs of RNA sequence and canonical base pairings taken from high-resolution structures. After the formation of native

secondary structures, we performed replica exchange DMD simulations with HRP-derived potentials. We then selected the top 100 structures based on low energy and high $C_{S-R}$ values and identified representative structures using a clustering analysis. Our goal was to define the RNA conformations that best represented the clusters (substates) of low-energy conformational ensembles that had strong correlations with the experimental HRP reactivities.
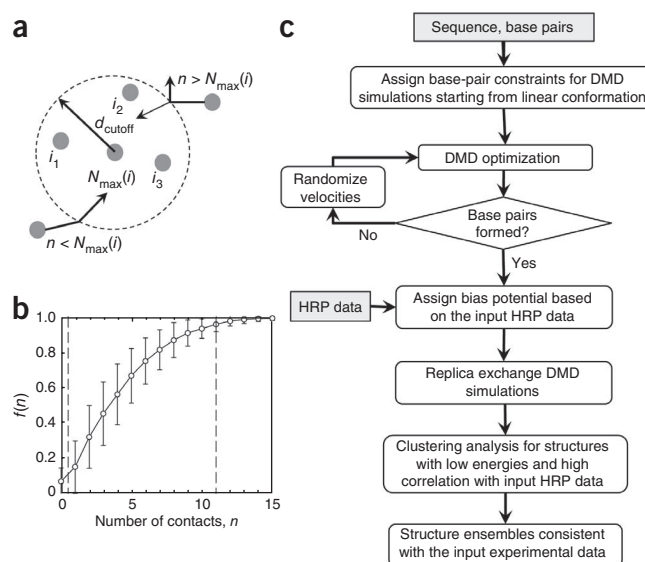
The training set for the initial DMD refinements were the yeast aspartic acid transfer (tRNA$^{Asp}$) (75 nt), the thiamine pyrophosphate (TPP) riboswitch (80 nt), the RNase P specificity domain (152 nt), the P546 domain of the *Tetrahymena thermophila* group I intron (158 nt), the M-box riboswitch (161 nt) and the *Azoarcus* group I intron (214 nt) RNAs. These six RNAs have diverse folds and different levels of higher-order packing interactions. The *Azoarcus* group I intron, M-box riboswitch and P546 domain RNAs have folds that are defined by close helical packing (**Fig. 3**); in contrast, the folds for the RNase P domain, the TPP riboswitch and tRNA$^{Asp}$ RNAs are characterized by local interactions between coaxially stacked helices (**Supplementary Fig. 2**). HRP is appropriate for *de novo* RNA structure refinement for the subset of RNAs with extensive helical packing. The extent of higher-order RNA packing can be estimated a priori from the fraction of nucleotides, $f(r)$, with HRP reactivities smaller than a given reactivity, $r$ (**Supplementary Fig. 3** and Online Methods). At $r = 0.25$, the RNAs with extensive helix packing interactions—the *Azoarcus* group I intron, M-box riboswitch and P546 domain RNAs—had significantly larger $f(r)$ values than the other RNAs (**Table 1** and **Supplementary Fig. 3**).

We characterized the predicted structural ensembles for each RNA in terms of the number and population of clusters in the 100 final structures. For each cluster, we also computed the root-mean-square deviation (RMSD) relative to the accepted structure and the prediction significance or P value[16] (**Table 1**). The RMSD value corresponding to a significant prediction varies with RNA size, and therefore it is not appropriate to apply a single cutoff for all RNAs[16]. For example, an RMSD of 7–10 Å is not significant for a small RNA but is highly significant for a 100-nt RNA[16]. The P value quantifies the statistical significance of the RNA fold prediction as the probability of observing a given conformation in an unbiased simulation with a pre-constrained base-pairing arrangement. $P < 0.01$ corresponds to predictions with high statistical significance[16].

We obtained highly significant predictions for each of the three largest RNAs in the training set. For the *Azoarcus* group I intron and the M-box riboswitch RNAs, all predicted structures
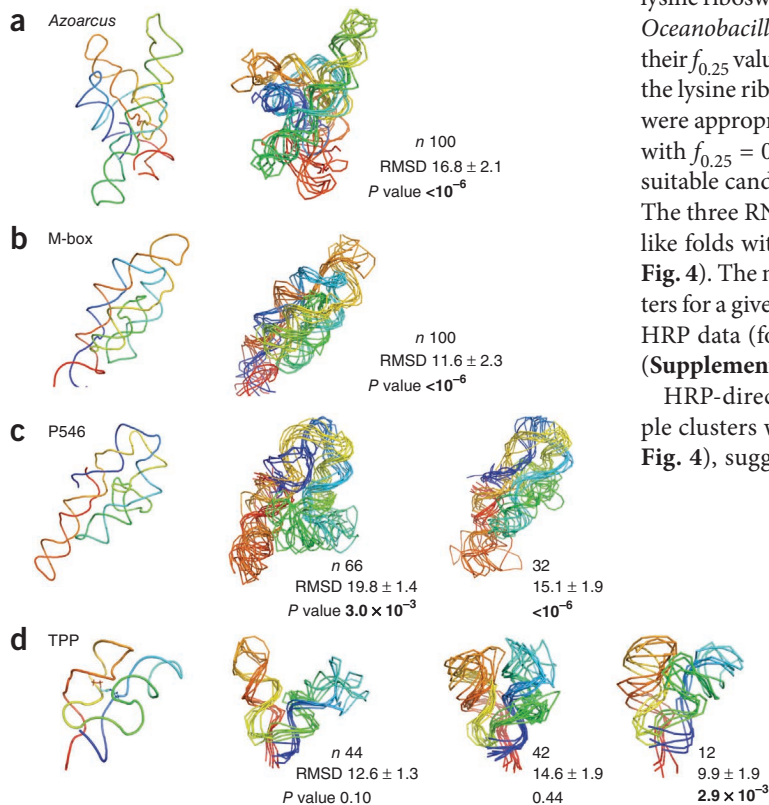
**Figure 2** | The assignment of potentials for incorporating HRP reactivities into DMD simulations. (**a**) Scheme for modeling the number of allowed contacts. Each nucleotide is assigned a threshold number of contacts ($N_{max}$) within the cutoff distance ($d_{cutoff} = 14$ Å). For a given nucleotide $i$, its $n$ through-space neighbors are denoted as $i_1$, $i_2$, $i_3$... An approaching nucleotide can form a new contact (indicated by the inward arrow) if the number of total contacts is smaller than $N_{max}$. If $n$ is larger than $N_{max}$, the approaching nucleotide can form a contact only if the total DMD kinetic energy is sufficient to overcome the energy penalty for overpacking (Online Methods). Otherwise, the nucleotide reflects back without forming a new contact (denoted by the outward arrow). (**b**) The fraction of nucleotides, $f(n)$, forming, at most, a given number of contacts, $n$. The mean (open circles) and s.d. (error bars) were computed over all single-chain RNA structures in the Protein Data Bank. Adjacent and same-helix nucleotide neighbors were excluded from the calculations of the number of contacts. The vertical dashed lines correspond to the minimal and maximal number of contacts of 0.5 and 11, respectively. (**c**) The HRP-directed DMD simulation algorithm.



fell into a single cluster with a low average RMSD and low $P$ value and were thus native-like (**Fig. 3a,b**). For the P546 domain RNA, the refined structures formed two highly populated clusters, both of which had low $P$ values and differed primarily in the location of a single helix (**Fig. 3c**). Simulations of the TPP riboswitch RNA yielded three clusters of structures. $P$ values for two of these clusters were high ($P > 0.01$), although the third cluster had a significant $P$ value (0.003) and correctly recapitulated the TPP ligand-binding pocket (**Fig. 3d**). For the tRNA$^{Asp}$ and RNase P RNAs, the HRP-directed structure refinement did not generate native-like structures ($P > 0.01$; **Table 1** and **Supplementary Fig. 2**).

### Consequences of $f_{0.25}$ and HRP data quality for refinement

For the six training RNAs, we observed a strong correlation between the fraction of nucleotides protected from HRP

cleavage, $f_{0.25}$, and the prediction significance (**Table 1**). The tRNA$^{Asp}$, RNase P and TPP riboswitch RNAs had $f_{0.25}$ values less than 0.25 and yielded inaccurate predictions, whereas we obtained statistically significant fold predictions for the *Azoarcus* group I intron, M-box riboswitch and P546 domain RNAs, which had higher $f_{0.25}$ values (**Fig. 3**). The $f_{0.25}$ values are calculated based on the HRP data alone without reference to the accepted structure. Thus, we conclude that the HRP-directed structure refinement is appropriate for RNAs with extensive close packing of helices, corresponding to $f_{0.25} > 0.25$.

We next applied HRP-directed structure refinement to the test set of four additional RNA molecules: the glmS ribozyme (152 nt), lysine riboswitch (174 nt), catalytic domain of RNase P (231 nt) and *Oceanobacillus iheyensis* group II intron (412 nt) RNAs. Based on their $f_{0.25}$ values (**Table 1**), three of these RNAs—the glmS ribozyme, the lysine riboswitch and the catalytic domain of RNase P RNAs— were appropriate candidates for structure refinement. In contrast, with $f_{0.25} = 0.21$, the *O. iheyensis* group II intron RNA was not a suitable candidate for refinement using HRP-derived constraints. The three RNAs with appropriate $f_{0.25}$ values all refined to native-like folds with significant $P$ values (**Table 1** and **Supplementary Fig. 4**). The major structural variations between the different clusters for a given RNA corresponded to regions without well-defined HRP data (for example, the 3′ end of RNase P catalytic domain) (**Supplementary Fig. 4** and **Supplementary Dataset**).

HRP-directed structure predictions often resulted in multiple clusters with distinct structures (**Fig. 3** and **Supplementary Fig. 4**), suggesting that not all experimental constraints can be



**a** *Azoarcus*

*n* 100
RMSD 16.8 ± 2.1
*P* value <10$^{-6}$

**b** M-box

*n* 100
RMSD 11.6 ± 2.3
*P* value <10$^{-6}$

**c** P546

*n* 66                           32
RMSD 19.8 ± 1.4        15.1 ± 1.9
*P* value **3.0 × 10$^{-3}$**    <10$^{-6}$

**d** TPP

*n* 44            42            12
RMSD 12.6 ± 1.3   14.6 ± 1.9   9.9 ± 1.9
*P* value 0.10    0.44         **2.9 × 10$^{-3}$**

**Figure 3** | HRP-directed RNA fold refinement for the training set. RNAs are shown with backbone traces. On the left is the accepted structure for each RNA. On the right are representative structures for each highly populated cluster. Small clusters (with one or two structures) are not shown. Backbones are colored from blue to red in the 5′ to 3′ direction. For each cluster, the number of structures, the mean RMSD and the $P$ value are shown. Significant $P$ values[16] are emphasized in bold.

satisfied simultaneously. Predictions that yielded multiple clusters reflected either the intrinsic structural heterogeneity of an RNA molecule or non-ideal experimental data. To explore the relationship between prediction accuracy and experimental HRP data quality, we generated idealized datasets by assuming perfect structure-reactivity correlations ($C_{S-R} = 1$) for the M-box riboswitch, P546 domain and TPP riboswitch RNAs (Online Methods and **Supplementary Table 1**). We generated additional simulated datasets with decreasing $C_{S-R}$ values by introducing random noise into the idealized data. Larger $C_{S-R}$ values generally yielded significant increases in the RNA prediction significance (**Supplementary Table 1**).

## DISCUSSION

HRP-directed structure refinement is unique among the RNA structure refinement methods, as its prediction quality is highest for larger and more complex RNA folds that have extensive helical packing and that have a substantial fraction of their nucleotides occluded from solvent (**Table 1**, **Fig. 3** and **Supplementary Fig. 4**). The HRP-directed fold prediction is also highly tolerant of the noise intrinsic to the RNA HRP experiment. Although the overall correlations between structure and HRP reactivity, as illustrated by the $C_{S-R}$ values, were modest (**Fig. 1b**), we obtained highly significant refinements because our algorithm reduces the effect of noise by identifying subsets of nucleotides with a high probability of being buried or exposed (**Fig. 1c**) and imposes strong energy terms on these nucleotides to drive RNA folding.

Previous RNA tertiary structure prediction studies have shown that a relatively small number of long-range constraints are often sufficient to reduce the allowable conformational space and make the prediction of diverse native-like structures possible[8,10]. In three of the RNAs studied here, the *Azoarcus* group I intron, the lysine riboswitch and the glmS ribozyme RNAs, we included long-range pseudoknot base-pairing constraints in the structure prediction. Even for these partially pre-constrained RNAs, the HRP-directed structural refinement improved the predictions beyond what is possible with the inclusion of the pseudoknot base-pairing constraints alone (Online Methods). One can thus use HRP-directed structural refinement in conjunction with other classes of information. Moreover, the correlation between the structure prediction accuracy and the quality of the input HRP data (**Supplementary Table 1**) suggests that if it becomes possible to improve the HRP approach or if experiments that better measure the solvent accessibility of RNA molecules are developed, it will be possible to refine RNA folds with an even higher level of accuracy.

The goal of our method is to reconstruct structural models for challenging RNA molecules that are not amendable to high-resolution methods. These structural models are especially useful for developing experimentally testable hypotheses and for guiding the exploration of structure-function relationships in RNA.

All software packages developed in this work for analyzing hydroxyl radical data and for predicting RNA structural models are available at http://troll.med.unc.edu/ifoldrna/HRP-1.0-openmpi.tgz.

## METHODS

Methods and any associated references are available in the online version of the paper.

*Note: Supplementary information is available in the online version of the paper.*

**AUTHOR CONTRIBUTIONS**
F.D., K.M.W. and N.V.D. conceived of and designed the computational and experimental procedures. C.A.L. performed and analyzed the HRP measurements. F.D. developed the computational methodology and performed the computational analysis. F.D., C.A.L., K.M.W. and N.V.D. wrote the manuscript.

1. Gesteland, R.F., Cech, T.R. & Atkins, J.F. eds. *The RNA World: The Nature of Modern RNA Suggests a Prebiotic RNA World*. (Cold Spring Harbor Lab Press, Plainview, NY, 2006).
2. Das, R. & Baker, D. Automated *de novo* prediction of native-like RNA tertiary structures. *Proc. Natl. Acad. Sci. USA* **104**, 14664–14669 (2007).
3. Ding, F. *et al. Ab initio* RNA folding by discrete molecular dynamics: from structure prediction to folding mechanisms. *RNA* **14**, 1164–1173 (2008).
4. Parisien, M. & Major, F. The MC-Fold and MC-Sym pipeline infers RNA structure from sequence data. *Nature* **452**, 51–55 (2008).
5. Cao, S. & Chen, S.J. Physics-based *de novo* prediction of RNA 3D structures. *J. Phys. Chem. B* **115**, 4216–4226 (2011).
6. Das, R. *et al.* Structural inference of native and partially folded RNA by high-throughput contact mapping. *Proc. Natl. Acad. Sci. USA* **105**, 4144–4149 (2008).
7. Yu, E.T., Hawkins, A., Eaton, J. & Fabris, D. MS3D structural elucidation of the HIV-1 packaging signal. *Proc. Natl. Acad. Sci. USA* **105**, 12248–12253 (2008).
8. Gherghe, C.M. *et al.* Native-like RNA tertiary structures using a sequence-encoded cleavage agent and refinement by discrete molecular dynamics. *J. Am. Chem. Soc.* **131**, 2541–2546 (2009).
9. Jonikas, M.A. *et al.* Coarse-grained modeling of large RNA molecules with knowledge-based potentials and structural filters. *RNA* **15**, 189–199 (2009).
10. Lavender, C.A., Ding, F., Dokholyan, N.V. & Weeks, K.M. Robust and generic RNA modeling using inferred constraints: a structure for the hepatitis C virus IRES pseudoknot domain. *Biochemistry* **49**, 4931–4933 (2010).
11. Yang, S., Parisien, M., Major, F. & Roux, B. RNA structure determination using SAXS data. *J. Phys. Chem. B* **114**, 10039–10048 (2010).
12. Michel, F. & Westhof, E. Modelling of the three-dimensional architecture of group I catalytic introns based on comparative sequence analysis. *J. Mol. Biol.* **216**, 585–610 (1990).
13. Gutell, R.R. *et al.* Identifying constraints on the higher-order structure of RNA: continued development and application of comparative sequence analysis methods. *Nucleic Acids Res.* **20**, 5785–5795 (1992).
14. Deigan, K.E., Li, T.W., Mathews, D.H. & Weeks, K.M. Accurate SHAPE-directed RNA structure determination. *Proc. Natl. Acad. Sci. USA* **106**, 97–102 (2009).
15. Weeks, K.M. Advances in RNA structure analysis by chemical probing. *Curr. Opin. Struct. Biol.* **20**, 295–304 (2010).
16. Hajdin, C.E., Ding, F., Dokholyan, N.V. & Weeks, K.M. On the significance of an RNA tertiary structure prediction. *RNA* **16**, 1340–1349 (2010).
17. Bailor, M.H., Mustoe, A.M., Brooks, C.L. III & Al-Hashimi, H.M. Topological constraints: using RNA secondary structure to model 3D conformation, folding pathways, and dynamic adaptation. *Curr. Opin. Struct. Biol.* **21**, 296–305 (2011).
18. Tullius, T.D. & Greenbaum, J.A. Mapping nucleic acid structure by hydroxyl radical cleavage. *Curr. Opin. Chem. Biol.* **9**, 127–134 (2005).
19. Cate, J.H. *et al.* Crystal structure of a group I ribozyme domain: principles of RNA packing. *Science* **273**, 1678–1685 (1996).

20. Pastor, N., Weinstein, H., Jamison, E. & Brenowitz, M. A detailed interpretation of OH radical footprints in a TBP-DNA complex reveals the role of dynamics in the mechanism of sequence-specific binding. *J. Mol. Biol.* **304**, 55–68 (2000).

21. Bergman, N.H. *et al.* The three-dimensional architecture of the class I ligase ribozyme. *RNA* **10**, 176–184 (2004).

22. Rangan, P., Masquida, B., Westhof, E. & Woodson, S.A. Assembly of core helices and rapid tertiary folding of a small bacterial group I ribozyme. *Proc. Natl. Acad. Sci. USA* **100**, 1574–1579 (2003).

23. Dokholyan, N.V., Buldyrev, S.V., Stanley, H.E. & Shakhnovich, E.I. Discrete molecular dynamics studies of the folding of a protein-like model. *Fold. Des.* **3**, 577–587 (1998).

24. Dann, C.E. III *et al.* Structure and mechanism of a metal-sensing regulatory RNA. *Cell* **130**, 878–892 (2007).

25. Balasubramanian, B., Pogozelski, W.K. & Tullius, T.D. DNA strand breaking by the hydroxyl radical is governed by the accessible surface areas of the hydrogen atoms of the DNA backbone. *Proc. Natl. Acad. Sci. USA* **95**, 9738–9743 (1998).

26. Berman, H.M. *et al.* The Protein Data Bank. *Nucleic Acids Res.* **28**, 235–242 (2000).

27. Westhof, E., Dumas, P. & Moras, D. Restrained refinement of 2 crystalline forms of yeast aspartic-acid and phenylalanine transfer-RNA crystals. *Acta Crystallogr. A* **44**, 112–123 (1988).

28. Serganov, A. *et al.* Structural basis for gene regulation by a thiamine pyrophosphate-sensing riboswitch. *Nature* **441**, 1167–1171 (2006).

29. Krasilnikov, A.S., Yang, X., Pan, T. & Mondragon, A. Crystal structure of the specificity domain of ribonuclease P. *Nature* **421**, 760–764 (2003).

30. Adams, P.L. *et al.* Crystal structure of a self-splicing group I intron with both exons. *Nature* **430**, 45–50 (2004).

31. Cochrane, J.C., Lipchock, S.V., Smith, K.D. & Strobel, S.A. Structural and chemical basis for glucosamine 6-phosphate binding and activation of the glmS ribozyme. *Biochemistry* **48**, 3239–3246 (2009).

32. Serganov, A., Huang, L. & Patel, D.J. Structural insights into amino acid binding and gene control by a lysine riboswitch. *Nature* **455**, 1263–1267 (2008).

33. Kazantsev, A.V., Krivenko, A.A. & Pace, N.R. Mapping metal-binding sites in the catalytic domain of bacterial RNase P RNA. *RNA* **15**, 266–276 (2009).

34. Toor, N. *et al.* Tertiary architecture of the *Oceanobacillus iheyensis* group II intron. *RNA* **16**, 57–69 (2010).

## ONLINE METHODS

**HRP measurements.** *RNA preparation.* RNAs were synthesized by T7 RNA polymerase-mediated *in vitro* transcription[35] using double-stranded PCR-generated templates. Sequences were transcribed in the context of 5′ and 3′ structure cassette sequences to facilitate analysis by primer extension[36]. Transcribed RNAs were purified on 10% denaturing polyacrylamide gels (7 M urea, 1 × Tris-borate-EDTA (TBE)). Bands containing full-length product were excised; RNA was recovered by passive elution in 1 × TE (10 mM Tris, pH 8.0, 1 mM EDTA) and precipitation with ethanol. Samples were resuspended in 1 × TE and quantified by absorbance measurements at 260 nm.

*Hydroxyl radical cleavage.* HRP datasets for the *Azoarcus* group I intron and the RNase P specificity domain RNAs were taken from a previous study, which used essentially the same approach as outlined below[37]. Hydroxyl radical cleavage experiments for the other RNAs were performed as described[24]. RNAs were first refolded by heat denaturation, snap cooling on ice and incubation at 37 °C. The HRP data reported here are consistent with previously reported experiments[24,29,38,39].

For the ligand-binding RNAs, 1 μl of a 5 μM RNA solution was combined with 2 μl sterile water and 3 μl folding mix (333 mM 4-(2-hydroxyethyl)-1-piperazineethanesulfonic acid (HEPES), pH 8.0, 333 mM NaCl and 33 mM $MgCl_2$ for the TPP riboswitch RNA; 333 mM HEPES, pH 8.0, 333 mM KCl and 33 mM $MgCl_2$ for the lysine riboswitch RNA; and 167 mM HEPES, pH 7.5 and 6.7 mM $MgCl_2$ for the glmS ribozyme RNA). RNAs were heated at 95 °C for 2 min, cooled on ice and then incubated at 37 °C for 10 min. One microliter of ligand solution (10 μM thiamine pyrophosphate, 50 μM lysine or 1 mM glucoasamine-6-phosphate for the TPP riboswitch, the lysine riboswitch and the glmS ribozyme RNA, respectively) was added, and the RNA was incubated in the presence of ligand at 37 °C for 20 min.

To fold the other RNAs, 1 μl of a 5 μM RNA solution was combined with 3 μl sterile water and 3 μl folding mix (46.6 mM HEPES, pH 7.5 and 23.3 mM $MgCl_2$ for the tRNA[Asp] RNA; 333 mM HEPES, pH 7.5, 333 mM NaCl and 33 mM $MgCl_2$ for the P546 domain RNA; 46.6 mM HEPES, pH 7.5 and 23.3 mM $MgCl_2$ for the M-box riboswitch RNA; 33 mM HEPES, pH 7.5, 333 mM NaCl and 33 mM $MgCl_2$ for the RNase P catalytic domain RNA; 333 mM HEPES, pH 8.0, 333 mM KCl and 416 mM $MgCl_2$ for the *Azoarcus* group I intron RNA; and 300 mM HEPES, pH 8.0, 300 mM KCl and 375 mM $MgCl_2$ for the *O. iheyensis* group II intron RNA). These RNAs were then heated at 95 °C for 2 min, cooled on ice and then incubated at 37 °C for 20 min.

The glmS ribozyme construct contained a point mutation (G40A) to prevent autolytic RNA cleavage during the HRP experiment; this mutant induces minimal structural disruption to the RNA[40].

Hydroxyl radical cleavage was initiated by adding Fe(II)-EDTA, sodium ascorbate and hydrogen peroxide to the folded RNA. Fresh Fe(II)-EDTA (10 mM ammonium Fe(II) sulfate and 20 mM EDTA, pH 8.0) and 50 mM sodium ascorbate solutions were made before each experiment. The Fe(II)-EDTA and ascorbate solutions were combined in a 1:1 ratio, and 2 μl of this 1:1 solution and 1 μl of 0.03% hydrogen peroxide were spotted in separate areas of the reaction lid. Hydroxyl radical cleavage was initiated by brief centrifugation. After incubation at 37 °C for 2 min, reactions were quenched by the addition of a solution containing 169 μl water,

20 μl 3 M sodium acetate (pH 5.5) and 1 μl 20 μg/μl glycogen, followed by the addition of 500 μl ethanol. Modified RNA was recovered by precipitation with ethanol and washed with 70% ethanol. For each reaction, a no-reaction control without Fe(II)-EDTA and ascorbate was performed in parallel.

*Primer extension.* Sites of hydroxyl-radical–mediated cleavages were analyzed by primer extension using fluorescently labeled primers[37,41] labeled with fluorophores from the Applied Biosystems G5 dye set: FAM was used for the positive (+) reaction; VIC was used for the negative (−) reaction; and NED was used as the sequencing ladder. For each primer extension reaction, 3 μl 0.3 μM fluorescently labeled DNA primer was added to 1 pmol RNA in 10 μl 0.5 × TE. This solution was incubated at 65 °C for 5 min and then cooled on ice for 1 min. To this solution, 6 μl Superscript III enzyme mix (250 mM KCl, 167 mM Tris, pH 8.3, 1.67 mM of each deoxynucleotide triphosphate, 17 mM dithiothreitol and 10 mM $MgCl_2$) and 1 μl Superscript III (Invitrogen) were added. For the sequencing reactions, 1.67 mM of a dideoxynucleotide triphosphate was included in the Superscript III enzyme mix. The solution was incubated at 45 °C for 1 min, at 52 °C for 25 min and at 65 °C for 5 min. The (+) and (−) reagent and sequencing reactions were then combined in 1 ml ethanol to quench extension and to precipitate the complementary DNA (cDNA). Recovered cDNAs were washed with 70% ethanol and resuspended in 10 μl dry formamide (Applied Biosystems).

cDNAs were resolved on Applied Biosystems 3130 Genetic Analyzer or 3500 Genetic Analyzer capillary electrophoresis instruments. Signal processing, sequencing alignment and peak integration of raw traces were performed using ShapeFinder[42] and custom signal-processing software. A representative processed electropherogram is provided in **Supplementary Figure 5**. The net reactivity at each nucleotide was defined as the area of the (+) reaction peak after subtracting the area of the corresponding (−) reaction peak. Nucleotides with high (−) signals were excluded from further analyses as high-background regions; the number of these high-background regions was small in the analyzed RNAs. Net reactivities were normalized by dividing the reactivities by the average reactivity of the top 10% of nucleotides, excluding the top 2%. HRP reactivities for each of the ten RNAs are provided in the **Supplementary Dataset**.

**Computational modeling using HRP reactivities.** *Overview of the DMD refinement approach.* We used a coarse-grained RNA model for the DMD simulations[3] in which each nucleotide was represented by three pseudoatoms representing the phosphate, sugar and base groups. Bonded terms, including covalent bond lengths, angles and dihedrals, were used to model the local RNA geometry. Nonbonded interactions included base pairing, base stacking, phosphate-phosphate repulsion and hydrophobic interactions. We explicitly modeled the entropy loss for loop formation. To bias the DMD simulation toward the structural ensemble that was consistent with the experimental measurements, we added additional potential terms based on the experimental hydroxyl-radical–probing data.

The DMD simulations and analyses were performed in three steps. First, serial DMD simulations were performed with inputs of RNA sequence and canonical base pairings, including pseudoknotted pairs, as obtained from high-resolution structures.

Although the base-pairing arrangements were taken from X-ray crystallographic analyses, this information can be obtained with high accuracy from a sequence covariation analysis[12,13] or with selective 2′-hydroxyl acylation analyzed by primer extension–directed secondary structure prediction[14,15]. The result of these simulations was the formation of native secondary structures. Second, replica exchange DMD simulations with the HRP-derived potentials were applied to enrich for conformations that were consistent with the experimental HRP data. Third, the top 100 structures were selected that had the lowest energies and highest correlations between the HRP reactivities and the number of contacts ($C_{\text{S-R}}$). A clustering analysis based on the pairwise RMSD values was then performed to identify representative structures that were consistent with the predicted structural ensemble.

*HRP bias potential.* For each nucleotide, we assigned a favorable energy increment, $E_{\text{attr}}(i)$, for forming a contact; a threshold number of contacts, $N_{\text{max}}(i)$; and a repulsive energy, $dE_{\text{rep}}(i)$, for exceeding the threshold. The HRP $E_{\text{bias}}$ potential equals

$$E_{\text{bias}} = \sum_{i<j} E_{ij}^{\text{attr}} + \sum_{i} E_{i}^{\text{overbury}} \tag{1}$$

The first term is the pairwise attraction, $E_{ij}^{\text{attr}}(r_{ij}) = \min\{E_{\text{attr}}(i), E_{\text{attr}}(j)\}F(r_{ij})$, where $F(x)$ is a step function

$$F(x) = \begin{cases} 0 & IR < r \\ 1 & R_{\text{hc}} < r \le IR \\ \infty & r \le R_{\text{hc}} \end{cases} \tag{2}$$

$IR$ is the interaction range of 14 Å, and $R_{\text{hc}}$ is the hard core diameter, 3.0 Å. The second term prevents overburying by exceeding the threshold number of contacts

$$E^{\text{overbury}}(i) = dE_{\text{rep}}(i)(n_c(i) - N_{\text{max}}(i))\Theta(n_c(i) - N_{\text{max}}(i)) \tag{3}$$

where $n_c(i)$ is the number of contacts for each nucleotide $i$, $dE_{\text{rep}}(i)$ is the penalty energy for overburying and $\Theta(x)$ is the unit step function, which equals 1 if $x$ is positive and zero otherwise. The number of contacts for each nucleotide was computed as the number of nonlocal sugar beads within the 14 Å cutoff. For each nucleotide $i$, we excluded contacts with nucleotides that were adjacent (within 4 nt) to $i$ or were adjacent to a nucleotide to which $i$ base pairs (for $i$ pairing with $I$, these nucleotides are $|i - j| > 4$ or $|I - j| > 4$).

*Assignment of interaction parameters.* The energy parameters, $N_{\text{max}}(i)$, $E_{\text{attr}}(i)$ and $dE_{\text{rep}}(i)$ were assigned using the HRP reactivities for each nucleotide in three steps. (i) Assignment of the threshold number of contacts. The threshold number of contacts, $N_{\text{max}}$, was assigned according to the reactivities, $R$, smoothed over a sliding window of 3 nt. Smoothing reduced the influence of the noise intrinsic to HRP experiments performed with RNA and increased the correlations to the accepted structure, $C_{\text{S-R}}$. We defined two threshold values, $R_{\text{min}}$ and $R_{\text{max}}$, corresponding to the maximally buried and exposed nucleotides, respectively. The $R_{\text{min}}$ and $R_{\text{max}}$ values were the average of the subsets from 2% to 20% and from 80% to 98%, respectively, of the rank-ordered $R$ values. The top and bottom 2% of the $R$ values were discarded

to reduce the influence of extreme $R$ values that are observed in typical HRP experiments. For nucleotides with $R$ values smaller than $R_{\text{min}}$ or higher than $R_{\text{max}}$, the threshold number of contacts was defined as $NC_{\text{max}} = 11$ and $NC_{\text{min}} = 0.5$, respectively (**Fig. 2b**). For a nucleotide $i$ with intermediate reactivity, the threshold number of contacts was assigned by linear interpolation

$$\begin{aligned}N_{\text{max}}(i) = &NC_{\text{max}} \\ &+ (NC_{\text{min}} - NC_{\text{max}})(R(i) - \langle R_{\text{min}}\rangle)/(\langle R_{\text{max}}\rangle - \langle R_{\text{min}}\rangle)\end{aligned} \tag{4}$$

(ii) Assignment of representative buried and exposed nucleotides. We first identified segments of strongly buried and exposed nucleotides. We defined three values, $R^{\text{EXP}}$, $R^{\text{INT}}$ and $R^{\text{BUR}}$, corresponding to the threshold values of exposed, intermediate and buried residues, respectively (**Fig. 1c**). The buried threshold $R^{\text{BUR}}$ and the exposed threshold $R^{\text{EXP}}$ correspond to the lowest 20% and the highest 20% of the rank-ordered reactivities, $R$, respectively. There are two types of intermediate $R$ values: the average value of all the reactivities, $\langle R\rangle$, and the median value of the rank-ordered reactivities, $R_{50}$. For simplicity, we chose the mean of these two values as $R^{\text{INT}}$.

We defined buried segments as those with more than three consecutive nucleotides having $R$ smaller than $R^{\text{INT}}$ and at least one nucleotide having $R$ smaller than $R^{\text{BUR}}$. For each buried segment, we selected the one with the lowest reactivity as the buried representative, excluding the first and last residues in the segment. Similarly, we defined exposed segments as those with more than three consecutive nucleotides having $R$ larger than $R^{\text{INT}}$ and at least one nucleotide having $R$ larger than $R^{\text{EXP}}$ and, for two-nucleotide segments, with both nucleotides having $R$ values larger than $R^{\text{EXP}}$. For each exposed segment, we defined the nucleotide with the largest $R$ value as the exposed representative.

(iii) Assignment of attractions and repulsions. Two attractive energy scales were used, $E_{\text{low}} = -0.10$ kcal/mol and $E_{\text{high}} = -0.05$ kcal/ml, based on the simulation temperature (see below). We assigned a strong attractive energy, $E_{\text{low}}$, to the buried representative nucleotides identified in step ii and the median value of $(E_{\text{high}} + E_{\text{low}})/2$ to their nearest neighbors. For all remaining nucleotides, we assigned the weak attractive energy of $E_{\text{high}}$. We defined a strong repulsive overburial energy, $dE_{\text{rep}}(i) = 0.3$ kcal/mol, for both the buried and exposed representative positions. We assigned the repulsive energy $dE_{\text{rep}}(i) = -E_{\text{attr}}(i)$ to all other nucleotides, where $E_{\text{attr}}(i)$ was equal to $E_{\text{low}}$ or $E_{\text{high}}$. By making the overburial repulsion potentials equal to those for the attractions, these nucleotides were allowed to make additional contacts ($>N_{\text{max}}$) without a net energy penalty. This approach reduced the effect of noise in the HRP experiments on RNA structure refinement by promoting compaction while imposing strong energy terms correlated with solvent accessibility for the subset of nucleotides identified as having a high probability of being buried or exposed. The HRP-derived values—threshold number of contacts ($N_{\text{max}}$), attractive energy ($E_{\text{attr}}$) and repulsive energy ($dE_{\text{rep}}(i)$)—are listed for all tested RNAs in the **Supplementary Dataset**.

*Replica exchange DMD simulations.* Because the HRP-directed potential is nonspecific with respect to any two nucleotides (in contrast to the distance and bonded constraints between specific nucleotides[8,10]), we performed replica exchange DMD simulations to obtain a sufficient sampling of conformational space. We used eight replicas with temperatures of 0.200, 0.225, 0.250,

0.270, 0.300, 0.333, 0.367 and 0.400 kcal/(mol·$k_B$). Every 1,000 DMD time units, we exchanged the replicas with neighboring temperatures according to a Metropolis-based Monte Carlo algorithm using instantaneous potential energies[3]. For each replica, we performed the simulations over $5 \times 10^5$ DMD time units. Replica exchange DMD simulations were performed in parallel on 2.27 GHz computing nodes (Intel Xeon). The representative running times for the TPP riboswitch (80 nt), M-box riboswitch (160 nt) and *Azoarcus* group I intron RNAs (214 nt) RNAs were 60, 170 and 264 CPU hours, respectively. The wall-clock time is one-eighth of the total CPU time.

*Identifying structure ensembles that were consistent with the experiments.* To identify structural ensembles that were consistent with the experimentally measured HRP data, we computed structure-reactivity correlation coefficients, $C_{S-R}$, for the snapshot structures, computed every 100 time units, yielding $4 \times 10^4$ snapshots for each refinement. We rank ordered these snapshots by their $C_{S-R}$ values and selected the 2,000 structures with the lowest (negative) correlation coefficients. From these, we selected the 100 structures with the lowest energies. We also selected structures by applying these rules in the reverse order: from the $4 \times 10^4$ total structures, we selected 2,000 structures with the lowest energy, from which we then selected the 100 structures with the lowest $C_{S-R}$ values.

For the combined 200 structures, we removed duplicates and selected the top 100 structures to represent the predicted structural ensemble. The structures were ranked according to their combined rank order using both their energy and $C_{S-R}$ values. We clustered these 100 structures according to their pairwise RMSDs using a hierarchical clustering algorithm and grouped similar structures into clusters using a cutoff RMSD. For simplicity, we used a cutoff value of 4 Å (~2 s.d.) below the average RMSD for a given RNA length[16] (see below), or three quarters of the average RMSD, whichever was smaller

$$\begin{cases} R(n) - 4, & R(n) - 4 < 0.75R(n) \\ 0.75R(n), & 0.75R(n) \le R(n) - 4 \end{cases} \quad (5)$$

Here, $n$ is the RNA length, and $R(n)$ is the average RMSD as the function of RNA length.

*P value calculations.* A recent study of a large set of RNA decoy structures derived from both simulations and threading suggested that the RMSD between two random RNA structures of the same length follows a Gaussian distribution with a length-dependent average RMSD and a length-independent s.d. (~1.8 Å)[16]. For an RNA with a known secondary structure, the average RMSD between two random decoy structures is smaller relative to a decoy set generated without knowledge of the secondary structure. We computed the statistical significance, or $P$ value, corresponding to the probability that an HRP-constrained structure prediction, evaluated by its RMSD from the accepted structure, is significantly better than that expected by chance. The $P$ value calculation tool is available online at http://ifoldrna.dokhlab.org[16].

The issue of how to interpret the significance of a structure model with a given RMSD value has been a major challenge in the field of RNA folding. Some groups have suggested that RMSDs should correspond in some qualitative way with the physical

dimensions of the RNA. For example, the RMSDs should be less than 7 Å (the average distance between two nucleotides) or within the width of an RNA helix (~20 Å). In fact, the average RMSD between any two structural models is strongly dependent on RNA length and whether the secondary structure is used as a constraint[16]. Thus, we argue that an appropriate way to understand the significance of a structure prediction is in terms of a $P$ value. Prior work using the 7 Å or 20 Å heuristic rules tended to overestimate the quality of the predictions for short RNAs and to underestimate the significance of the predictions for large RNAs. For large RNAs, seemingly large RMSD values with low $P$ values correspond to native-like folds with high significance (**Fig. 3**).

*Generation of ideal and randomized reactivity profiles.* We generated idealized HRP reactivities based on the number of contacts in the native structure, $R^{ideal}(i) = 1 - N_c(i)/N_{max}$. We added noise to these idealized reactivities to generate randomized reactivities, $R^{rand}(i) = R^{ideal}(i)(1 + \sigma x)$, where $x$ is a random number from $-1$ to 1, and $\sigma$ is the amplitude of the noise, determined by the relative error

$$1/N \sum_{i=1}^{N} |(R^{rand}(i) - R^{ideal}(i))/R^{ideal}(i)| = 1/N \sum_{i=1}^{N} \sigma |x| \sim \sigma/2 \quad (6)$$

where the sum is over all the nucleotides in an RNA. By varying $\sigma$, we generated randomized reactivity profiles with different levels of noise and, thus, different structure-reactivity correlations (**Supplementary Fig. 6a**). Notably, the M-box riboswitch RNA had the least noise-induced decrease in the structure-reactivity correlation, $C_{S-R}$, whereas the tRNA$^{Asp}$ RNA had the greatest decrease in $C_{S-R}$ value, which correlates with their respective prediction significances (**Supplementary Fig. 6b**).

For the M-box riboswitch, P546 domain, TPP riboswitch and tRNA$^{Asp}$ RNAs, we selected seven sets of computationally generated HRP data with $C_{S-R}$ values ranging from $-0.4$ to $-1.0$ (**Supplementary Table 1**). Using the generated HRP reactivities as the input, we applied our structure refinement protocol to generate structural ensembles (**Supplementary Table 1**). For all the tested RNAs, except the tRNA$^{Asp}$ RNA, we found that the HRP reactivities with high $C_{S-R}$ values resulted in low RMSDs and highly significant predictions. As the $C_{S-R}$ values of the input HRP reactivities decreased, the RMSDs of the predicted structures and the corresponding $P$ values increased, indicating less accurate predictions.

There are two major implications of this analysis. First, the high $P$ value predictions for the tRNA$^{Asp}$ RNA using both experimental and computationally generated HRP reactivities suggest that RNAs, like tRNAs, with few buried nucleotides are not good candidates for HRP-directed refinement. Notably, these RNAs can be identified (and excluded) in advance using the $f_{0.25}$ metric (**Supplementary Fig. 3**). Second, our simulations indicate that the level of noise and the resulting structure-reactivity correlation for the input HRP data have a determining role in the accuracy of the HRP-directed structure prediction. If a better experimental method with reduced noise in HRP (or solvent accessibility) reactivities were developed, our approach would immediately lead to statistically significantly more accurate RNA structure refinements.

*Structural refinement for RNAs with pseudoknot base pairs.* In our study, we assumed that all base pairs, including pseudoknots,

were known. A relatively small number of constraints based on long-range contacts, such as pseudoknots, are sufficient to direct the prediction of highly significant RNA structures[10]. The *Azoarcus* group I intron, the lysine riboswitch and the glmS ribozyme RNAs contain long-range pseudoknots that probably reduce the available conformational space and may themselves lead to significant structure predictions. To examine the effects of incorporating HRP data into RNA refinements in which long-range pseudoknot constraints were included, we compared the results of RNA structure prediction with and without the HRP data.

First, we evaluated whether incorporation of the HRP data as constraints drives the conformational sampling toward native states during the course of the simulations for the pseudoknot-containing RNAs. We calculated the RMSDs for all RNA conformations sampled during the DMD simulations both with and without the HRP data as constraints. For both the lysine riboswitch and glmS ribozyme RNAs, incorporation of the HRP data into the DMD simulations significantly enhanced sampling of native-like conformations (**Supplementary Fig. 7**). Second, we applied the structure selection approach to reconstruct conformational ensembles for simulations that did not incorporate the HRP data. For these large RNAs, if the HRP data were not used to drive the refinement, the resulting structural ensembles fell into multiple small clusters with a wide range of RMSD values (**Supplementary Table 2**); in contrast, using the HRP data to drive the refinement yielded only a few clusters, each with well-defined structures and highly significant RMSD values (**Table 1**). Therefore, although the pseudoknotted base pairs reduced the available conformational space, the HRP-directed structural refinement drove RNA folding to native-like states.

35. Milligan, J.F., Groebe, D.R., Witherell, G.W. & Uhlenbeck, O.C. Oligoribonucleotide synthesis using T7 RNA polymerase and synthetic DNA templates. *Nucleic Acids Res.* **15**, 8783–8798 (1987).
36. Merino, E.J., Wilkinson, K.A., Coughlan, J.L. & Weeks, K.M. RNA structure analysis at single nucleotide resolution by Selective 2'-Hydroxyl Acylation and Primer Extension (SHAPE). *J. Am. Chem. Soc.* **127**, 4223–4231 (2005).
37. Duncan, C.D.S. & Weeks, K.M. The Mrs1 splicing factor binds the bI3 group I intron at each of two tetraloop-receptor motifs. *PLoS One* **5**, e8983 (2010).
38. Murphy, F.L. & Cech, T.R. An independently folding domain of RNA tertiary structure within the *Tetrahymena* ribozyme. *Biochemistry* **32**, 5291–5300 (1993).
39. Latham, J.A. & Cech, T.R. Defining the inside and outside of a catalytic RNA molecule. *Science* **245**, 276–282 (1989).
40. Klein, D.J., Been, M.D. & Ferre-D'Amare, A.R. Essential role of an active-site guanine in glmS ribozyme catalysis. *J. Am. Chem. Soc.* **129**, 14858–14859 (2007).
41. McGinnis, J.L., Duncan, C.D. & Weeks, K.M. High-throughput SHAPE and hydroxyl radical analysis of RNA structure and ribonucleoprotein assembly. *Methods Enzymol.* **468**, 67–89 (2009).
42. Vasa, S.M. *et al.* ShapeFinder: a software system for high-throughput quantitative analysis of nucleic acid reactivity information resolved by capillary electrophoresis. *RNA* **14**, 1979–1990 (2008).