

Article

Statistical analysis of SHAPE-directed RNA secondary structure modeling

Srinivas Ramachandran, Feng Ding, Kevin M. Weeks, and Nikolay V. Dokholyan

Biochemistry, **Just Accepted Manuscript** • DOI: 10.1021/bi300756s • Publication Date (Web): 04 Jan 2013

Downloaded from <http://pubs.acs.org> on January 4, 2013

Just Accepted

“Just Accepted” manuscripts have been peer-reviewed and accepted for publication. They are posted online prior to technical editing, formatting for publication and author proofing. The American Chemical Society provides “Just Accepted” as a free service to the research community to expedite the dissemination of scientific material as soon as possible after acceptance. “Just Accepted” manuscripts appear in full in PDF format accompanied by an HTML abstract. “Just Accepted” manuscripts have been fully peer reviewed, but should not be considered the official version of record. They are accessible to all readers and citable by the Digital Object Identifier (DOI®). “Just Accepted” is an optional service offered to authors. Therefore, the “Just Accepted” Web site may not include all articles that will be published in the journal. After a manuscript is technically edited and formatted, it will be removed from the “Just Accepted” Web site and published as an ASAP article. Note that technical editing may introduce minor changes to the manuscript text and/or graphics which could affect content, and all legal disclaimers and ethical guidelines that apply to the journal pertain. ACS cannot be held responsible for errors or consequences arising from the use of information contained in these “Just Accepted” manuscripts.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Statistical analysis of SHAPE-directed RNA secondary structure modeling

Srinivas Ramachandran^{1,‡}, Feng Ding^{1, 3,‡}, Kevin M. Weeks² and Nikolay V. Dokholyan^{1,}*

¹Department of Biochemistry and Biophysics, ²Department of Chemistry, University of North Carolina at Chapel Hill, ³Department of Physics and Astronomy, Clemson University

*Corresponding Author: Nikolay V. Dokholyan, 3097 Genetic Medicine Building, Campus Box 7260, Chapel Hill, NC 27599. E-mail: dokh@unc.edu, Phone: 919-843-2513.

‡These authors contributed equally.

Funding: This work was supported by grants from the NIH (GM080742 to N.V.D. and AI068462 to K.M.W.).

1
2
3
4
5 **ABSTRACT:** The ability to predict RNA secondary structure is fundamental for understanding
6
7 and manipulating RNA function. The structural information obtained from selective 2'-hydroxyl
8
9 acylation analyzed by primer extension (SHAPE) experiments greatly improves the accuracy of
10
11 RNA secondary structure prediction. Recently, Das and colleagues [Kladwang *et al.*,
12
13 *Biochemistry* **50**:8049 (2011)] proposed a “bootstrapping” approach to estimate the variance and
14
15 helix-by-helix confidence levels of predicted secondary structures based on resampling
16
17 (randomizing and summing) the measured SHAPE data. We show that the specific resampling
18
19 approach described by Kladwang *et al.* introduces systematic errors and underestimates
20
21 confidence in secondary structure prediction using SHAPE data. Instead, a leave-data-out
22
23 jackknife approach better estimates the influence of a given experimental dataset on SHAPE-
24
25 directed secondary structure modeling. Even when 35% of the data were left out in the jackknife
26
27 approach, the confidence levels of SHAPE-directed secondary structure prediction were
28
29 significantly higher than those calculated by Das and colleagues using bootstrapping. Helix
30
31 confidence levels were thus significantly underestimated in the recent study, and resampling
32
33 approach implemented by Kladwang *et al.* is not an appropriate metric for assigning confidences
34
35 in SHAPE-directed secondary structure modeling.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 Despite an explosion in discoveries of RNAs and their functional roles in biology, accurate
4 knowledge of structures of these molecules is incomplete.¹ Knowledge of information encoded in
5 RNA structures, especially in RNA secondary structures (the pattern of base pairs) is necessary
6 for understanding and manipulating RNA function. Computational RNA secondary structure
7 prediction methods^{2,3} have been widely used to generate structural hypotheses in RNA research.
8 Secondary structures predicted from sequence alone often have significant errors, however,
9 including both falsely predicted and missing base pairs.^{1,4,5} Incorporation of experimental
10 structural information derived from chemical probing experiments can significantly improve
11 secondary structure predictions.⁶ For example, the comprehensive and quantitative information
12 available from selective 2'-hydroxyl acylation analyzed by primer extension (SHAPE) probing
13 experiments greatly improves the accuracy of RNA secondary structure prediction.^{5,7,8}

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29 With widespread adoption of SHAPE and other experimental approaches for directing RNA
30 secondary structure prediction, rigorous *a priori* estimation of the confidence level of a given
31 RNA secondary structure prediction would constitute a major and welcome advance. Recently,
32 Das and colleagues proposed a “bootstrapping” (resampling) approach to estimate the variance
33 and confidence level of predicted secondary structures based on resampling of measured SHAPE
34 data.⁹ Based on their statistical study, Das and colleagues suggested that the confidence level of
35 SHAPE-derived RNA secondary structure prediction is about 77%. Follow up analysis of the
36 work of Das and colleagues revealed that important components of their experimental work were
37 not consistent with recommended practices in using and evaluating SHAPE technologies.¹⁰ In
38 this correspondence, we show that the specific resampling approach developed by Das and
39 colleagues is unphysical, introduces systematic error into the resampled data, and results in a
40 large underestimation of the confidence of SHAPE-directed secondary structure prediction. As
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 detailed here, a leave-data-out jackknife approach more accurately estimates the influence of a
4 given experimental dataset on SHAPE-directed secondary structure modeling.
5
6
7
8
9

10 MATERIALS AND METHODS

11
12 **SHAPE-directed secondary structure prediction.** We compared the resampling approach of
13 Kladwang et al. and our jackknife approach on four RNA molecules: tRNA^{Phe} (76 nucleotides),
14 adenine riboswitch (71 nucleotides), cyclic-di-GMP riboswitch (97 nucleotides), and 5S rRNA
15 (120 nucleotides). The SHAPE data for these RNAs are presented in the companion manuscript¹⁰
16 and were used to direct secondary structure prediction using RNAstructure as described.⁵
17
18
19
20
21
22
23

24 **Kladwang resampling.** For RNA molecules with nucleotide positions 1, 2, 3,..., N , bootstrap
25 decoys were generated as described⁹. To summarize, N nucleotide positions were randomly
26 picked with repetition. Thus, some nucleotide positions were picked multiple times whereas
27 others were not picked at all. If a nucleotide position was picked three times, for example, the
28 SHAPE reactivity of that position was multiplied by 3. In total, 400 such decoy SHAPE datasets
29 were generated and the RNA secondary structure predicted for each decoy SHAPE dataset using
30 RNAstructure.
31
32
33
34
35
36
37
38
39

40 **Jackknife.** The jackknife decoy was generated by picking $(1 - f) \times N$ nucleotide positions
41 randomly where f is the fraction of data omitted. We performed jackknife analysis using f values
42 of 0.1 and 0.35. Thus, for a 100-nucleotide RNA molecule, the jackknife decoy dataset will
43 contain SHAPE values for 90 and 65 nucleotides, respectively, picked randomly out of the 100
44 values. Four hundred such decoy datasets were generated and the RNA secondary structure
45 predicted for each decoy SHAPE dataset using RNAstructure. For both jackknife and
46 bootstrapping approaches, we calculated the percentage of these 400 decoy RNA secondary
47 structures that contained each base-pair observed in the secondary structure predicted by the
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 original SHAPE data. The generation of a decoy dataset for each approach is illustrated in Figure
4
5
6 1. Control simulations performed with 100 decoy datasets converged with those using 400
7
8 decoys; we report the results of the 400 member decoy datasets here. Performing a jackknife
9
10 analysis with 35% omitted data represents an extreme case, because it is uncommon to miss so
11
12 many SHAPE measurements. Thus, the jackknifing results presented here represent an over-
13
14 estimation of the variability of SHAPE-directed structure modeling.
15
16
17
18
19

20 RESULTS AND DISCUSSION

21
22 SHAPE reactivity at a position (S_i) in an RNA reports the conformational flexibility of a given
23
24 nucleotide and is inversely related to the propensity of the position to form a base pair or tertiary
25
26 contact.^{7,11} The measured SHAPE reactivity can be used as an experimentally-based correction to
27
28 bias an RNA secondary structure prediction.⁵ To estimate confidences for individual helices in a
29
30 given secondary structure prediction, Das and colleagues proposed a bootstrapping approach to
31
32 resample the measured SHAPE data.⁹ This approach was motivated by the demonstrated utility
33
34 of bootstrapping analyses in evaluation of phylogenetic trees.¹² Bootstrapping entails reshuffling
35
36 a given set of data with repetition; critically, for bootstrapping to be statistically justified the
37
38 calculated quantity usually needs to be independent of the order of the shuffled data, an
39
40 assumption that is generally valid in the construction of phylogenetic trees.
41
42
43
44

45
46 In the approach taken by Das and colleagues, SHAPE reactivities were generated for an RNA
47
48 molecule by “resampling with repetition” the original SHAPE reactivities. Given SHAPE
49
50 reactivities S_i for nucleotides $i = 1, 2, 3, \dots, N$, resampling with repetition entails picking N
51
52 indices randomly, whereby an index value can be selected multiple times for a given position,
53
54 such that some nucleotides are not picked at all, whereas some are picked multiple times. This
55
56 specific implementation is not appropriate for treatment of SHAPE data as SHAPE reactivities
57
58
59
60

1
2
3 are *not* independent: *both* the actual reactivity and the position associated with a given reactivity
4 are required for secondary structure modeling. As SHAPE data are thus inherently unsuitable for
5
6 “resampling with replacement”, Das and colleagues retained the positional information of
7
8 SHAPE reactivity after resampling (hence, there was no shuffling *per se*).⁹ In addition, when a
9
10 position was sampled multiple times during bootstrapping, the SHAPE reactivity of the position
11
12 was increased by the amount equal to the original reactivity of the position each time (Das,
13
14 personal communication). Hence, the relative error introduced at each position in this resampling
15
16 approach is $n \times S_i$, where n is the number of repetitions, (0, 1, 2,...; Fig. 1A, B). This treatment of
17
18 repetition is intrinsically different from the repetition introduced by a standard bootstrapping
19
20 approach as implemented for multiple sequence alignments, which increases the weight of the
21
22 repeated column in the alignment and thus enhances *both* sequence similarity and differences. In
23
24 contrast, in the approach used by Das and colleagues, repetition of high SHAPE values increases
25
26 the probability of breaking a base pair, whereas repetition of low SHAPE data does not increase
27
28 the probability of a forming a base pair.

29
30
31
32
33
34
35
36 This problem is most significant in regions of low, but non-zero, SHAPE reactivities that are
37
38 often base paired and was compounded by SHAPE signal processing errors in the prior work.¹⁰
39
40 Several occurrences of multiplying the SHAPE reactivity value under this approach would
41
42 effectively destabilize a helix in which the original SHAPE data would otherwise clearly score as
43
44 unreactive, and likely paired, overall. Thus, the bootstrapping approach as implemented by Das
45
46 and colleagues⁹ results in introduction of systematic error because the errors are integer multiples
47
48 of the original data itself. In addition, the decoy datasets generated were not the same size
49
50 because, after each round of shuffling to retain positional information, the number of data points
51
52 differs from the number in the original dataset. Since the errors added to the dataset are integer
53
54 multiples of the data itself, the perturbation is high, resulting in many datasets that are highly
55
56
57
58
59
60

1
2
3 dissimilar to the experimental data. Overall, this approach results in an underestimation of the
4 confidence of secondary structure prediction from the SHAPE data.
5
6

7
8 An alternate statistical method for estimation of the influence of a specific experimental
9 dataset on secondary structure modeling is the jackknife approach. Resampling by the jackknife
10 approach has also been used to estimate the confidence of phylogenetic trees.¹³ With jackknifing,
11 the fluctuation is introduced simply through omission of a fraction of the data, and the number of
12 data points in each decoy dataset is identical. The key variable in jackknife resampling is the
13 fraction of excluded data.
14
15
16
17
18
19
20
21

22 To quantify the extent of underestimation of the confidence of SHAPE prediction by the
23 method used by Das and colleagues, we used the jackknife method to evaluate recovery of
24 SHAPE-predicted base pairs for four RNA molecules – tRNA^{Phe}, the adenine and cyclic di-GMP
25 riboswitches, and 5S rRNA by generating decoys with up to 35% of the SHAPE data randomly
26 removed. SHAPE data were measured using the documented approach for performing the
27 SHAPE experiment,^{8,14} which is substantially different from that introduced by Das and
28 colleagues.^{9,10} Experimentally, it is uncommon to have 35% of SHAPE missing and few
29 experimentalists would try to predict a structure with such a high level of missing information,
30 we therefore also performed a jackknife analysis with 10% of the data omitted, the latter
31 corresponding to a more realistic practical occurrence.
32
33
34
35
36
37
38
39
40
41
42
43
44

45 In general, both the Kladwang et al. resampling approach and jackknifing identified the same
46 helices as being less well-defined by the SHAPE data (Fig. 2). However, the Kladwang et al.
47 approach grossly underestimated the confidence of the helices in three of the SHAPE-directed
48 structure models, those of tRNA^{Phe}, the cyclic-di-GMP riboswitch, and 5S rRNA. Estimated
49 confidences based on bootstrapping were less than confidence estimates obtained using decoy
50 datasets that were missing (an extreme) 35% of the experimental data (Fig. 2). In contrast, using
51
52
53
54
55
56
57
58
59
60

1
2
3 the jackknife approach with a more physically and experimentally realistic 10% omitted data, the
4 majority of the base pairs have 100% prediction confidences, and the lowest is ~80%, supporting
5
6 the general robustness of the SHAPE-directed secondary structure models.
7
8

9
10 A second problem with the specific resampling approach created by Das and coworkers lies
11 in their helix-by-helix interpretation. The bootstrap resampling of the SHAPE data introduces
12 noise that perturbs the relative free energy of each RNA structure. Longer helices with lower free
13 energies are less sensitive to this perturbation and are more likely to have high confidence in
14 prediction, especially given that Das and colleagues define the bootstrap value of a helix as the
15 maximum of the bootstrap values of base pairs across that helix.⁹ By this definition, longer
16 helices have more base pairs and are more likely to have at least one base pair with a high
17 “bootstrap value”. In contrast, shorter and less stable helices are more sensitive to perturbations
18 and thus more prone to break under perturbation. Therefore, the estimated bootstrap value of a
19 helix primarily reflects the stability of the helix rather than the underlying SHAPE data.
20
21
22
23
24
25
26
27
28
29
30
31
32

33 In summary, the specific bootstrapping procedure proposed by Das and colleagues to
34 resample the SHAPE data is unphysical, introduces systematic error into the resampled data, and
35 results in an underestimation of the confidence of SHAPE-directed secondary structure modeling.
36 The calculated confidence level obtained via this approach is not an appropriate metric for
37 estimation of the accuracy of experimentally-directed RNA secondary structure prediction.
38 Instead, this work supports use of the jackknife approach to generate resampled SHAPE data and
39 to estimate the sensitivity of predicted secondary structures to the underlying SHAPE dataset.
40 The more general issue of *a priori* identification of individual highly probable helices within a
41 given experimentally-directed RNA structure model remains a major research challenge.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

REFERENCES

1. Shapiro, B. A., Yingling, Y. G., Kasprzak, W., and Bindewald, E. (2007) Bridging the gap in RNA structure prediction, *Curr Opin Struct Biol* 17, 157-165.
2. Mathews, D. H. (2006) Revolutions in RNA secondary structure prediction, *J Mol Biol* 359, 526-532.
3. Mathews, D. H., Sabina, J., Zuker, M., and Turner, D. H. (1999) Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure, *J Mol Biol* 288, 911-940.
4. Mathews, D. H., and Turner, D. H. (2006) Prediction of RNA secondary structure by free energy minimization, *Curr Opin Struct Biol* 16, 270-278.
5. Deigan, K. E., Li, T. W., Mathews, D. H., and Weeks, K. M. (2009) Accurate SHAPE-directed RNA structure determination, *Proc. Natl. Acad. Sci. USA* 106, 97-102.
6. Mathews, D. H., Disney, M. D., Childs, J. L., Schroeder, S. J., Zuker, M., and Turner, D. H. (2004) Incorporating chemical modification constraints into a dynamic programming algorithm for prediction of RNA secondary structure, *Proc. Natl. Acad. Sci. USA* 101, 7287-7292.
7. Merino, E. J., Wilkinson, K. A., Coughlan, J. L., and Weeks, K. M. (2005) RNA structure analysis at single nucleotide resolution by Selective 2'-Hydroxyl Acylation and Primer Extension (SHAPE), *J. Am. Chem. Soc.* 127, 4223-4231.

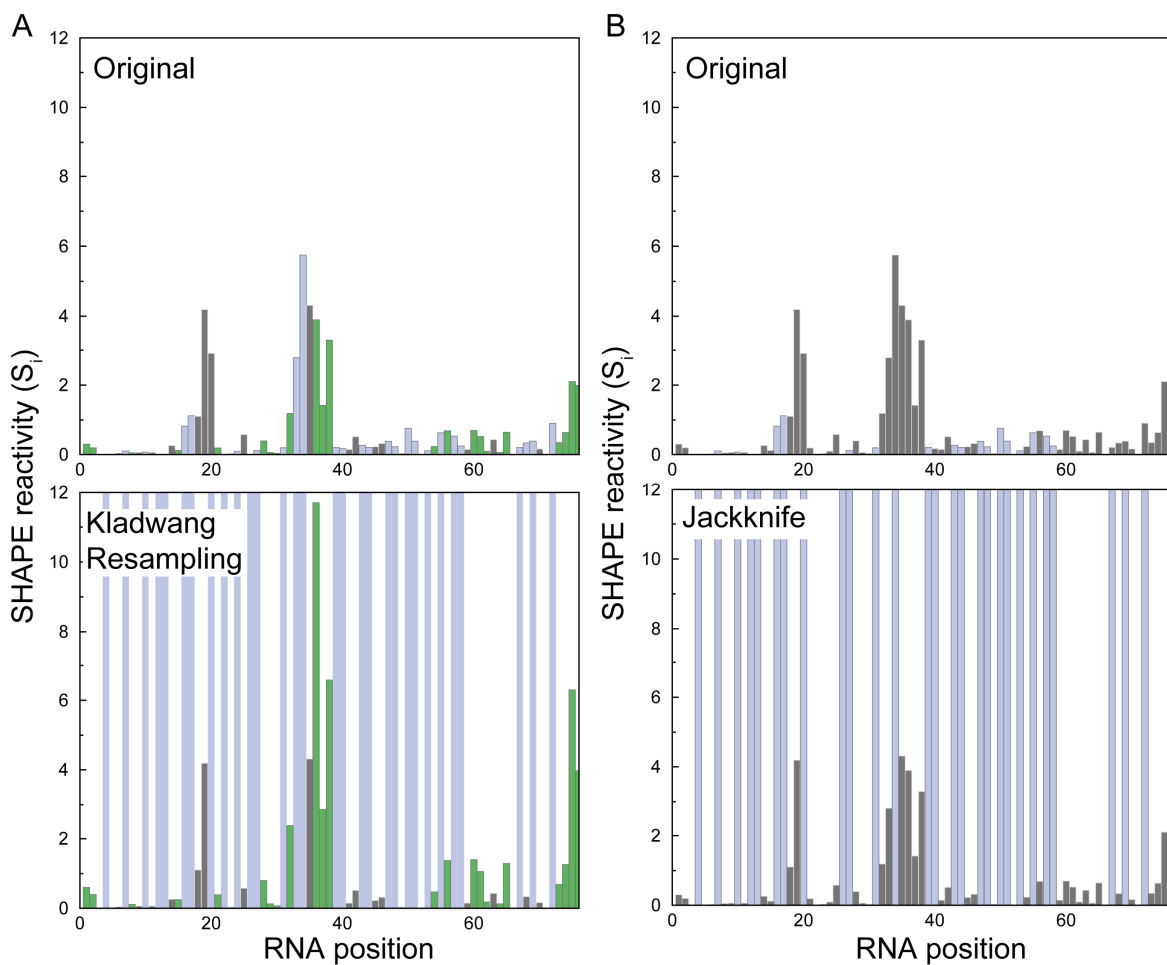
- 1
2
3 8. Wilkinson, K. A., Merino, E. J., and Weeks, K. M. (2006) Selective 2'-hydroxyl acylation
4 analyzed by primer extension (SHAPE): quantitative RNA structure analysis at single nucleotide
5 resolution, *Nat. Protocols 1*, 1610-1616.
6
7
- 8
9
10
11 9. Kladwang, W., VanLang, C. C., Cordero, P., and Das, R. (2011) Understanding the errors
12 of SHAPE-directed RNA structure modeling, *Biochemistry 50*, 8049-8056.
13
14
- 15
16
17 10. Leonard, C. W., Hajdin, C. E., Karabiber, F., Mathews, D. H., Favorov, O. V.,
18 Dokholyan, N. V. and Weeks, K. M. (2013) Principles for understanding the accuracy of
19 SHAPE-directed RNA structure modeling. *Biochemistry 52*, dx.doi.org/10.1021/bi300755u.
20
21
22
- 23
24
25 11. Gherghe, C. M., Shajani, Z., Wilkinson, K. A., Varani, G., and Weeks, K. M. (2008)
26 Strong correlation between SHAPE chemistry and the generalized NMR order parameter (S^2) in
27 RNA, *J Am Chem Soc 130*, 12244-12245.
28
29
30
- 31
32
33 12. Efron, B., Halloran, E., and Holmes, S. (1996) Bootstrap confidence levels for
34 phylogenetic trees, *Proc Natl Acad Sci U S A 93*, 7085-7090.
35
36
37
- 38
39
40 13. Soltis, P. S., and Soltis, D. E. (2003) Applying the bootstrap in phylogeny reconstruction,
41 *Statistical Science 18*, 256-267.
42
43
- 44
45
46 14. Vasa, S. M., Guex, N., Wilkinson, K. A., Weeks, K. M., and Giddings, M. C. (2008)
47 ShapeFinder: a software system for high-throughput quantitative analysis of nucleic acid
48 reactivity information resolved by capillary electrophoresis, *RNA 14*, 1979-1990.
49
50
51
52
53
54
55
56
57
58
59
60

Figure Legends

Figure 1. Generation of decoy SHAPE datasets via resampling-bootstrap and jackknife approaches. (A) tRNA^{Phe} SHAPE reactivities determined using a standard¹⁰ approach (top) and representative decoy dataset generated by specific resampling algorithm of Kladwang et al.⁹ (bottom). The SHAPE reactivities not present in the decoy datasets (top) are shown as blue bars in the original dataset and are shaded blue in the decoy dataset (bottom). Kladwang resampling also results in modified SHAPE reactivities (green bars, notice the increase in SHAPE reactivities at these positions in the decoy dataset). (B) SHAPE reactivities (top) and representative decoy dataset as generated by jackknifing (bottom). The jackknife approach results in SHAPE reactivities that are not picked in the decoy dataset (blue bars).

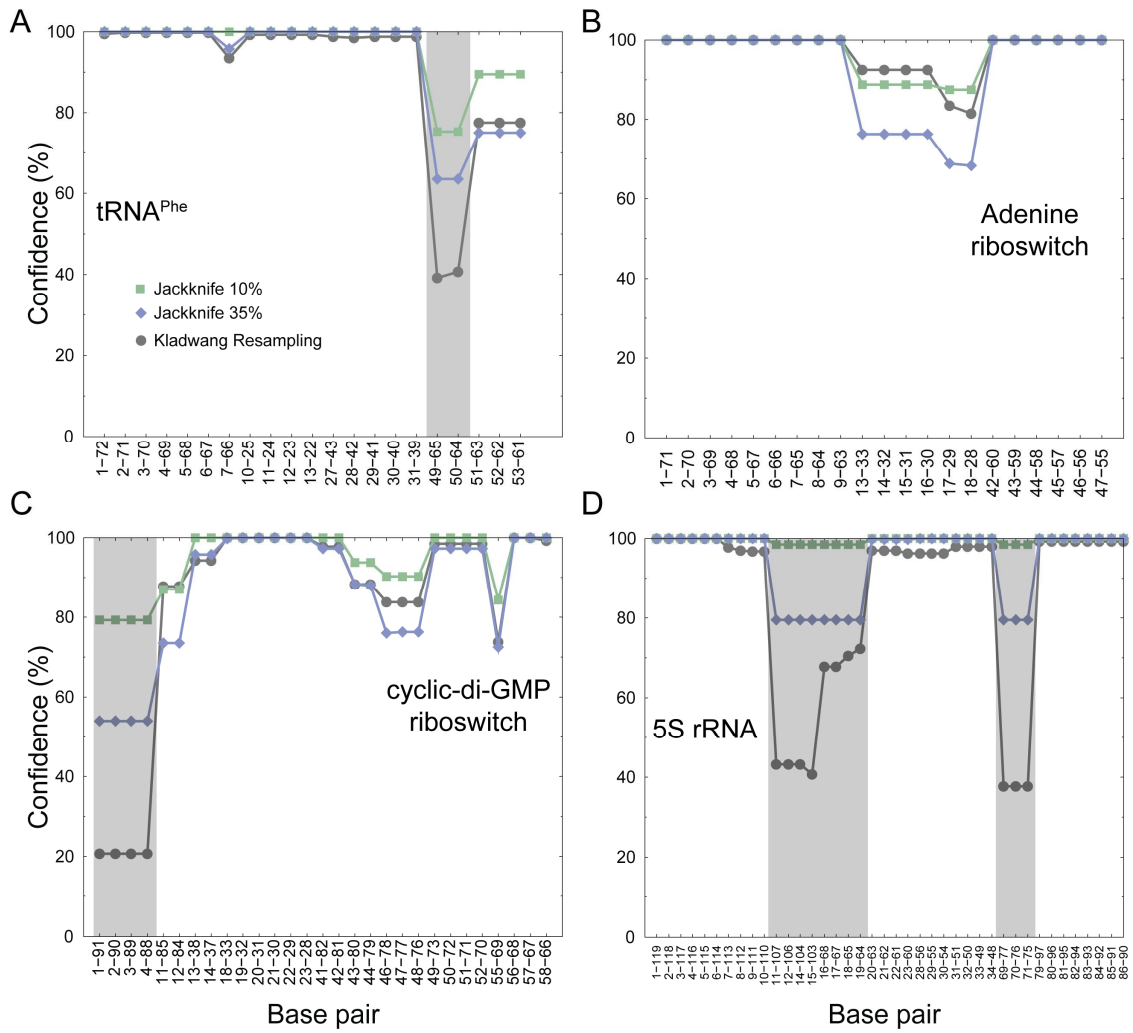
Figure 2. Confidence estimation for SHAPE-directed secondary structure modeling. Confidence estimates for SHAPE-predicted base pairs, calculated using the specific resampling algorithm of Kladwang et al.⁹ (gray circles) and using jackknife procedures that omit either 10% (green squares) or 35% (blue diamonds) of the data for (A) tRNA^{Phe}, (B) adenine riboswitch, (C) cyclic-di-GMP riboswitch, (D) 5S rRNA. Gray shading emphasizes regions in which the resampling-bootstrap approach underestimated confidences as compared to removing (an extreme) 35% of the experimental SHAPE data.

Figure 1



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 2



TOC Graphic

Statistical analysis of SHAPE-directed RNA secondary structure modeling

Srinivas Ramachandran, Feng Ding, Kevin M. Weeks and Nikolay V. Dokholyan

