1

2

# Covariation analysis with improved parameters reveals conservation in lncRNA structures

5

**Rafael C A Tavares[1], Anna Marie Pyle[1-3] and Srinivas Somarowthu[4]**

[1]Department of Chemistry, Yale University, New Haven, Connecticut, USA. [2]Department of Molecular, Cellular and Developmental Biology, Yale University, New Haven, Connecticut, USA. [3]Howard Hughes Medical Institute, Chevy Chase, Maryland, USA. [4]Department of Biochemistry and Molecular Biology, Drexel University College of Medicine, Philadelphia, Pennsylvania, USA.

11

12

13

14

15

16

17

18

19

## Abstract

The existence of phylogenetic covariation in base-pairing is strong evidence for functional elements of RNA structure, although available tools for identifying covariation are limited. R-scape is a recently developed program for prediction of covariation from sequence alignments, but it has limited utility on long RNAs, especially those of eukaryotic origin. Here we show that R-scape can be adapted for powerful prediction of covariation in long RNA molecules, including mammalian lncRNAs.

## Main

28    Long non-coding RNAs (lncRNAs) are well accepted as crucial regulators of gene

29    expression and disease progression[1]. Despite the ubiquity and significance of lncRNAs, our

30    understanding of structure-function relationships within this class of molecules is extremely

31    limited[2]. Studies of ribozymes, riboswitches, viral RNAs, mRNA UTRs and even coding

32    sequences have shown that conserved RNA secondary and tertiary structures are vital for RNA

33    function[3, 4]. It has therefore been of interest to determine whether lncRNA molecules contain

34    regions of functional structure and whether these structures are conserved[5-7]. If conservation in

35    base-pairing could be established, it would provide powerful evidence that RNA structure plays a

36    role in aspects of lncRNA function. Several empirical studies have demonstrated the existence of

37    structured regions within lncRNAs, and conventional phylogenetic covariation analyses were

38    found to support the empirically-determined structures[8-10]. Indeed in at least two cases, these

39    modules of RNA structure were flanked by highly conserved sequences that are consistent with a

40    biological role for lncRNA substructures[9, 10].

41    However, a powerful new method for stringent determination of nucleotide covariation,

42    known as R-scape, failed to support the existence of conserved base-pairings in well-studied

43    functional lncRNAs such as Xist and HOTAIR[11]. On the basis of these findings, it was

44    concluded that those lncRNAs do not contain conserved structure and are therefore unlikely to

45    contain functional elements of discrete structure. Like many tools for phylogenetic analysis, R-

46    scape was developed for application to small, highly structured RNA molecules for which many

47    sequences are available (such as bacterial riboswitches). We reasoned that, at least in its current

48    form, R-scape might not be equipped to confront the challenges posed by large, multidomain

49    eukaryotic RNA molecules. We therefore set out to test the limitations of R-scape covariation

3

50    analysis and to determine whether the approach could actually be modified in order to

51    successfully identify conservation of structures in mammalian lncRNAs.

52      A major challenge for the analysis of eukaryotic lncRNAs is the severe limitation in

53    available sequences[12]. We reasoned that this limitation, rather than any inherent lack of evidence

54    for lncRNA structure, might explain the reported inability of R-scape to identify conserved

55    structure in mammalian lncRNAs. To test this hypothesis, we analyzed the ability of R-scape to

56    detect basepair covariation in seven well-characterized, highly structured RNA molecules (tRNA,

57    5S ribosomal RNA, 5.8S ribosomal RNA, eukaryotic RNase P, U2 snRNA, U5 snRNA and the

58    eukaryotic small subunit ribosomal RNA) using input alignments that were restricted in three

59    different ways:  1) Inclusion of the original RFAM seed alignment 2) Sub-sampled alignments

60    and 3) Restriction to mammalian sequences. In the sub-sampled RFAM alignments, we limited

61    the number of sequences and the average pairwise identity to control for effects arising solely

62    from restrictions in these parameters (See Methods, Figure 1). The alignments restricted to

63    mammalian sequences represent the currently available alignments that have been built for most

64    lncRNAs.  Not surprisingly, there is a precipitous drop in covariation support for most of these

65    test RNAs in both the 'sub-sampled' and 'mammalian sequence' conditions.  Eukaryotic RNase

66    P (> 300 nt) is the most dramatic example, as only 13% of the base pairs can be flagged as

67    covariant by R-scape in the sub-sampling analysis (Figure 1). It is also worth highlighting the

68    particular case of 5.8S rRNA, for which the RFAM seed alignment already has a relatively high

69    pairwise sequence identity (~68%).  Predictably, R-scape finds covariation support for only 44%

70    of the base pairs in the 5.8S rRNA structure, and no support (0%) upon restriction of the analysis

71    to mammalian sequences. In fact, with the exception of tRNA, for which even mammalian

72    sequences have high nucleotide diversity, R-scape was unable to detect the majority of covarying

<div align="center">4</div>

73    base pairings in these model RNAs when the input alignments were limited to mammals. These

74    results indicate that R-scape fails to detect covariation not just in lncRNAs, but in most of the

75    structurally complex, well-characterized functional RNA molecules that have been tested.

76        It is important to note that RFAM alignments are hand-curated and refined[13], therefore,

77    deviations from RFAM's ideal heuristics may bias R-scape results. This phenomenon was shown

78    to be true for other covariance prediction algorithms when RFAM alignments were compared to

79    emulated genomic alignments as inputs[14]. Multiple sequence-based alignments from datasets like

80    the TBA/Multiz (UCSC genome browser) can be used to build covariation models and generate

81    structural alignments for lncRNAs, but these alignments lack the quality of RFAM alignments,

82    which can then affect R-scape prediction sensitivity. Finally, since genomic alignments may not

83    accurately reflect the regions of lncRNA loci that are actively expressed, there is a consistent

84    need for direct characterization and annotation of lncRNA transcripts across species in order to

85    improve identification of conserved sequence and structure motifs, as described elsewhere[15, 16].

86        There is accumulating evidence that lncRNAs possess local modules of RNA structure

87    and that they can contain both structured and unstructured regions[8, 10]. Given that R-scape uses

88    the entire length of an RNA sequence for analysis, it is possible that the presence of unstructured

89    regions negatively impacts the ability of R-scape to identify structural conservation. To test this,

90    we analyzed the ability of R-scape to predict covariation when unstructured regions are included

91    in an alignment. R-scape is reported to perform well on riboswitches, using sequences that are

92    restricted to the functional, structured region of the molecule. We therefore chose the SAM-I

93    riboswitch (RF00162) as an example, but we now included the surrounding mRNA regions from

94    the alignment. The mRNA regions were aligned using MAFFT[17], and the alignment for the

95    SAM-I riboswitch region was kept the same as in the RFAM alignment. We then compared R-

96    scape predictions by varying the number of sequences in the alignment. In the case of the SAM-I

97    riboswitch alone, R-scape predicted significant covariation even with only 40 sequences in the

98    alignment (Figure 2A), as reported previously. However, inclusion of the flanking mRNA in the

99    alignment resulted in a notable decrease in R-scape performance:  Even when 60 sequences are

100   included in the alignment, R-scape could identify covarying base pairs in only one helix (Figure

101   2B), indicating that the presence of unstructured RNA regions has a strong influence on R-scape

102   analysis output. However, as the number of sequences in the alignment increases, R-scape can

103   identify more covarying base pairs, even when unstructured regions are included. This suggests

104   that R-scape may ultimately become a powerful tool for identifying covarying base-pairs when a

105   sufficient number of sequences are provided ($> 90$ for SAM-I riboswitch). However, since the

106   alignments for most human lncRNAs are currently limited to 30-60 mammalian sequences, R-

107   scape default settings should not be applied to lncRNA covariation analysis.

108         Another feature that is expected to influence the performance of any covariation analysis

109   is the length of an RNA molecule and of its corresponding structural alignment.  LncRNAs are

110   typically very large and many exceed 1kb[18]. However, R-scape was benchmarked with a test set

111   consisting predominantly of small RNAs.  Of the 104 RNAs in that test set[11], there are only 21

112   RNAs with an average length greater than 200 nts and only seven that exceed 1kb, and all the

113   seven are ribosomal RNAs.  It is therefore unlikely that the R-scape default parameters are

114   appropriate for analysis of large RNAs. To test this, we asked whether R-scape performs better

115   when the analysis is broken down in short overlapping windows tiling the entire RNA rather than

116   when given a long whole-length alignment. We examined alignments (see methods) of two long

117   RNAs in sliding windows: 1) 7SK RNA (RFAM ID:00100) and 2) Aphthovirus internal

118   ribosome entry site (RFAM ID: 00210). For both RNAs, R-scape was able to identify more

119    covarying base-pairs when the analysis was run with sliding windows than when given the full-

120    length alignment (Fig. 3), indicating that the R-scape default parameters work better on short

121    alignments, either as aligned sequences of inherently small RNAs or long RNA alignments that

122    have been analyzed in a set of sliding windows.

123         Taken together (Figures 1-3), these results suggest that one might be able to increase the

124    signal-to-noise ratio for predicting lncRNA covariations by maximizing the number of sequences

125    (increasing alignment depth) and running R-scape analysis in short windows. Here, we applied

126    both conditions to analyze the RepA region of lncRNA Xist. In a previous study, R-scape

127    identified no significant base pair covariation in RepA structure[11]. However, the input alignment

128    in that study was limited to ten sequences, which was beneath an empirical threshold value (~40

129    sequences) suggested in the very same paper. We therefore reanalyzed RepA using a recent,

130    experimentally determined secondary structure[10] and we included significantly more sequences

131    in the alignment. As expected, just by adding more sequences we were able to identify

132    covariation in RepA, but it was limited to a single base pair. Interestingly, this base pair is

133    located within the functionally important repeat-five region[19]. To further improve the signal-to-

134    noise ratio, we ran R-scape on short (500-nt) overlapping windows, tiling the entire RNA. Using

135    this procedure, R-scape identified five statistically significant covariant base pairs: two in

136    domain I and three in domain II of the lncRNA RepA (Figure 4). Importantly, each of the

137    covariant base pairs detected are in the regions of high structural confidence (low Shannon

138    entropy, see Liu et al.,[10] and Supplementary figure 2). It is also worth highlighting that the three

139    base pairs in domain II identified by R-scape are in proximity to long-range crosslink sites

140    identified by Liu et al. 2017[10], and to a stretch of conserved base sequence, suggesting that even

7

141    though R-scape identified only five base pairs, they are consistent with experimental studies and

142    are likely to be functionally important.

143        Up to this point, our analysis suggests that the default parameters in R-scape are

144    exceedingly stringent and that they may not be sufficiently sensitive to predict covariation with

145    reduced alignment depth and low phylogenetic diversity, which are features inherent to most

146    current lncRNA alignments (Xist, HOTAIR, SRA, etc). Most telling, R-scape failed to detect

147    significant covariation when faced with similar alignments even for well-structured RNAs such

148    as ribosomal RNAs, snRNAs and the eukaryotic ribozyme RNAseP, suggesting that more

149    sequencing data is required to provide sufficient alignment depth for lncRNA structural

150    conservation analysis on R-scape. Given the plethora of lncRNA genes and their implicated roles

151    in human diseases, there is an urgent need for better tools and metrics to identify conserved

152    structures and associated functions of these giant molecules.

153        We therefore asked whether other metrics could improve the performance of R-scape on

154    long RNA molecules. RNAalifold with stacking[20] ($B^s_{i,j}$ renamed in Rivas et al., as RAFS) was

155    previously shown to be among the best performing covariation metric available and it has been

156    extensively validated in several RNA structure prediction platforms, where it is frequently

157    combined with structural stability metrics[21-23]. However, Rivas et al[11] have argued that the G-test

158    statistic (GT) performs better than RAFS in terms of positive predictive value (PPV) and thus

159    would be less prone to false positive discovery. To get a sense of the tradeoff between sensitivity

160    and PPV within these two metrics, we reanalyzed the original R-scape test set (104 RFAM

161    alignments) with default parameters. We used average product correction (APC) as it was shown

162    to improve the performance of both GT and RAFS (both renamed, then, as APC-GT and APC-

163    RAFS)[11]. First, we measured the difference in sensitivity and PPV of these two metrics by

8

164    varying the E-value threshold (Supplementary Figure 3). The PPV value for APC-RAFS gets

165    worse than APC-GT (> 5%) only for relatively high E-values (> 0.1).  However, at the default E-

166    value threshold of 0.05, APC-RAFS resulted in much higher sensitivity (~84%) relative to APC-

167    GT (~64%), with a PPV compromise of less than 4%, suggesting that APC-RAFS is in fact a

168    more robust metric than APC-GT.

169        Next, we tested the performance of these two metrics by varying the number of

170    sequences in the input alignment (Supplementary Figure 3). Most strikingly, APC-RAFS

171    achieved 63% sensitivity with only 20 sequences in the alignment compared to APC-GT, which

172    resulted in only 40% sensitivity with the same input. We then used APC-RAFS to score the same

173    alignments from Figure 1 (Supplementary Figs. 4 and 5) and observed a significant improvement

174    in covariation detection under restricted conditions (fewer sequences, increased average pairwise

175    identity and decreased phylogenetic diversity), relative to the original analysis using APC-GT.

176    Remarkably, the eukaryotic RNAseP case showed a dramatic 45% sensitivity increase upon

177    subsampled alignment analysis with APC-RAFS relative to APC-GT, and an even higher

178    improvement (49%) on the mammalian sequence alignment. In all cases, the use of APC-RAFS

179    on restricted alignments improved the overall covariation output when compared to APC-GT

180    with no compromise to specificity as given by PPV (Supplementary Fig. 5), indicating that APC-

181    RAFS is able to at least partly overcome the negative effects of lncRNA-like restrictions on R-

182    scape predictive power while preserving statistical rigor. All these observations suggest that

183    APC-RAFS is a highly robust metric for RNA covariation analysis with R-scape (null-model

184    based analysis) and, most importantly, the most suitable method for alignments with the

185    restrictions normally found in lncRNAs.

186        Based on the above observations, we utilized APC-RAFS to analyze the published

187    structural alignments for full-length lncRNA-RepA and Domain I of lncRNA-HOTAIR[9] (Fig. 5)

188    and found that R-scape is now able to support covariation of numerous base pairs in both RNAs.

189    We identified 16 covariant base pairs within the full-length lncRNA-RepA when the alignment

190    was analyzed in overlapping 500-nt windows tiling the RNA every 100 nt.  In this case, 9 out of

191    10 helical motifs with covariant base pairs flagged by R-scape/APC-RAFS were also suggested

192    to be conserved in previous empirical studies[10].   Within HOTAIR domain I, 24 base pairs were

193    flagged as covariant by R-scape/APC-RAFS in 10 helical segments of this region. Also, in this

194    case, most helices where APC-RAFS found covariant base pairs overlapped with helices

195    previously suggested as structurally conserved in domain I of HOTAIR[9]. These results strongly

196    suggest that APC-RAFS can be used within R-scape to improve covariation analysis of lncRNA

197    structure, confirming the conclusions from previous studies and highlighting the presence of

198    conserved structured regions in lncRNAs HOTAIR and RepA.

199        In conclusion, we show that R-scape default parameters are not applicable to lncRNAs,

200    but that R-scape is capable of identifying covariation when appropriately parameterized. We

201    suggest that increased alignment depth, sliding windows approach and a more sensitive statistical

202    metric, the APC-RAFS, are parameters that may help R-scape to identify conserved structural

203    elements in large molecules such as lncRNAs. By combining these approaches, we were able to

204    detect significant covarying base pairs in the experimental structures of lncRNAs HOTAIR and

205    RepA.  We hope that the results and approaches reported here provide improved tools for

206    meeting the challenges inherent to studying lncRNA molecules and that they facilitate future

207    studies and method development.

208

## Methods

### R-scape analysis

Seed alignments for tRNA, 5S ribosomal RNA, 5.8S ribosomal RNA, eukaryotic RNase P, U2 snRNA, U5 snRNA, small subunit ribosomal RNA (SS rRNA), 7SK, Aphthovirus IRES and SAM-I Riboswitch were downloaded from the RFAM database (RFAM v13.0). To obtain alignments restricted to mammals, mammalian sequences were manually extracted from each RNA family in the Rfam database and then aligned using Infernal (version 1.1.2). Sub-sampling analysis was performed by randomly selecting sequences using the 'submsa' option. The average pairwise identity (figure 1) was controlled using the 'maxid' option. The parameters and RFAM family IDs for all original and derived alignments are listed in Supplementary Figure 1. All analyses using R-scape were carried out at the default E-value (0.05), unless otherwise specified in the text.

Sliding window analyses were carried out with the "window" and "slide" options on R-scape, to define window size and sliding step of the R-scape search, respectively. Window size was varied between 50 - 500 nt, depending on the RNA length and structure, thereby ensuring that intact helices could be contained within the chosen window size.

The original R-scape test set was downloaded from the Eddy lab website (http://eddylab.org/R-scape/). The average product corrected RNAalifold with stacking (APC-RAFS) and G-test (APC-GT) statistics were compared using the "RAFSp" and "GTp" options respectively, by varying E-value thresholds and the number of sequences in the alignment.

11

## Acknowledgments

We thank Dr. Thayne Dickey for thoughtful comments on the manuscript. S.S. is supported by start-up funds from Drexel University College of Medicine and a CURE grant from the Pennsylvania Department of Health. R.C.A.T. was supported by NIH Grant RO1 50313. A.M.P. is an investigator with the Howard Hughes Medical Institute.

## Author contributions

S.S., R.C.A.T and A.M.P designed research and wrote the manuscript. S.S. and R.C.A.T carried out the experiments.

## Competing financial interests

The authors declare no competing financial interests.

# References

1. Schmitt, A.M. & Chang, H.Y. Long Noncoding RNAs: At the Intersection of Cancer and Chromatin Biology. *Cold Spring Harbor Perspectives in Medicine* **7** (2017).

2. Pyle, Anna M. Looking at LncRNAs with the Ribozyme Toolkit. *Molecular Cell* **56**, 13-17 (2014).

3. Mustoe, A.M. et al. Pervasive Regulatory Functions of mRNA Structure Revealed by High-Resolution SHAPE Probing. *Cell* **173**, 181-195 e118 (2018).

4. Pirakitikulr, N. et al. The Coding Region of the HCV Genome Contains a Network of Regulatory RNA Structures. *Mol Cell* **62**, 111-120 (2016).

5. Fu, Y. et al. Discovery of Novel ncRNA Sequences in Multiple Genome Alignments on the Basis of Conserved and Stable Secondary Structures. *PloS one* **10**, e0130200 (2015).

6. Will, S. et al. LocARNA-P: accurate boundary prediction and improved detection of structural RNAs. *Rna* **18**, 900-914 (2012).

7. Gruber, A.R. et al. RNAz 2.0: improved noncoding RNA detection. *Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing*, 69-79 (2010).

8. Novikova, I.V., Hennelly, S.P. & Sanbonmatsu, K.Y. Structural architecture of the human long non-coding RNA, steroid receptor RNA activator. *Nucleic Acids Research* **40**, 5034-5051 (2012).

9. Somarowthu, S. et al. HOTAIR Forms an Intricate and Modular Secondary Structure. *Molecular Cell* **58**, 353-361 (2015).

10. Liu, F., Somarowthu, S. & Pyle, A.M. Visualizing the secondary and tertiary architectural domains of lncRNA RepA. *Nature Chemical Biology* **13**, 282 (2017).

11. Rivas, E., Clements, J. & Eddy, S.R. A statistical test for conserved RNA structure shows lack of evidence for structure in lncRNAs. *Nature Methods* **14**, 45 (2016).

12. Nitsche, A. & Stadler, P.F. Evolutionary clues in lncRNAs. *Wiley Interdisciplinary Reviews: RNA* **8**, e1376 (2017).

13. Kalvari, I. et al. Rfam 13.0: shifting to a genome-centric resource for non-coding RNA families. *Nucleic Acids Research* **46**, D335-D342 (2018).

14. Smith, M.A. et al. Widespread purifying selection on RNA structure in mammals. *Nucleic Acids Research* **41**, 8220-8236 (2013).

15. Hezroni, H. et al. Principles of long noncoding RNA evolution derived from direct comparison of transcriptomes in 17 species. *Cell Rep* **11**, 1110-1122 (2015).

16. Chillon, I. & Pyle, A.M. Inverted repeat Alu elements in the human lincRNA-p21 adopt a conserved secondary structure that regulates RNA function. *Nucleic Acids Res* **44**, 9462-9471 (2016).

17. Kuraku, S. et al. aLeaves facilitates on-demand exploration of metazoan gene family trees on MAFFT sequence alignment server with enhanced interactivity. *Nucleic Acids Research* **41**, W22-W28 (2013).

18. Palazzo, A.F. & Lee, E.S. Non-coding RNA: what is functional and what is junk? *Frontiers in Genetics* **6** (2015).

19. Pintacuda, G., Young, A.N. & Cerase, A. Function by Structure: Spotlights on Xist Long Non-coding RNA. *Frontiers in Molecular Biosciences* **4**, 90 (2017).

20. Lindgreen, S., Gardner, P.P. & Krogh, A. Measuring covariation in RNA alignments: physical realism improves information measures. *Bioinformatics* **22**, 2988-2995 (2006).

21. Washietl, S., Hofacker, I.L. & Stadler, P.F. Fast and reliable prediction of noncoding RNAs. *Proceedings of the National Academy of Sciences of the United States of America* **102**, 2454-2459 (2005).

287    22.    Hofacker, I.L. RNA consensus structure prediction with RNAalifold. *Methods in molecular biology*
288           **395**, 527-544 (2007).
289    23.    Bernhart, S.H. et al. RNAalifold: improved consensus structure prediction for RNA alignments.
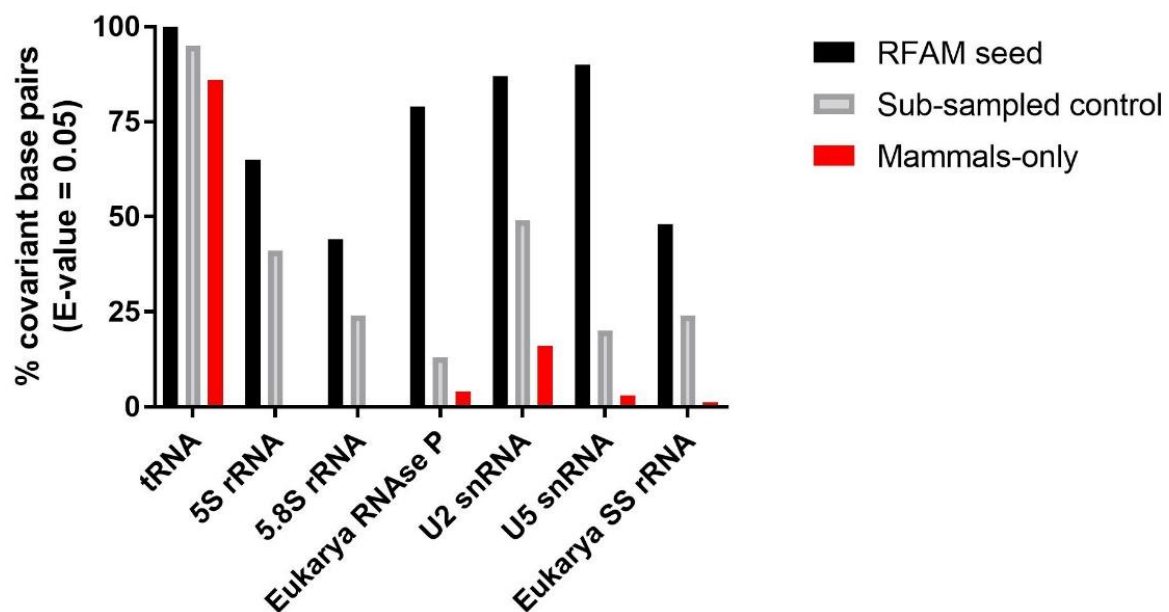290           *BMC bioinformatics* **9**, 474 (2008).
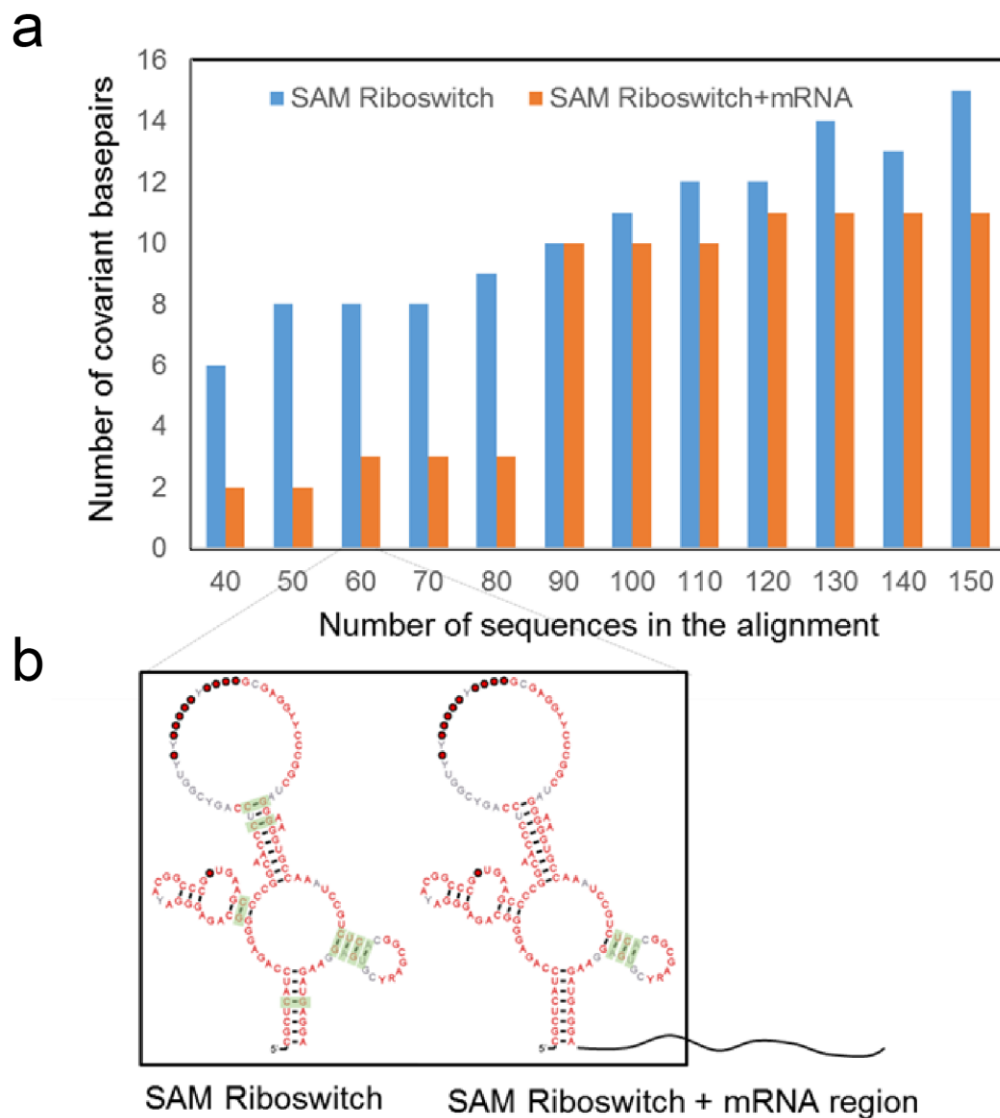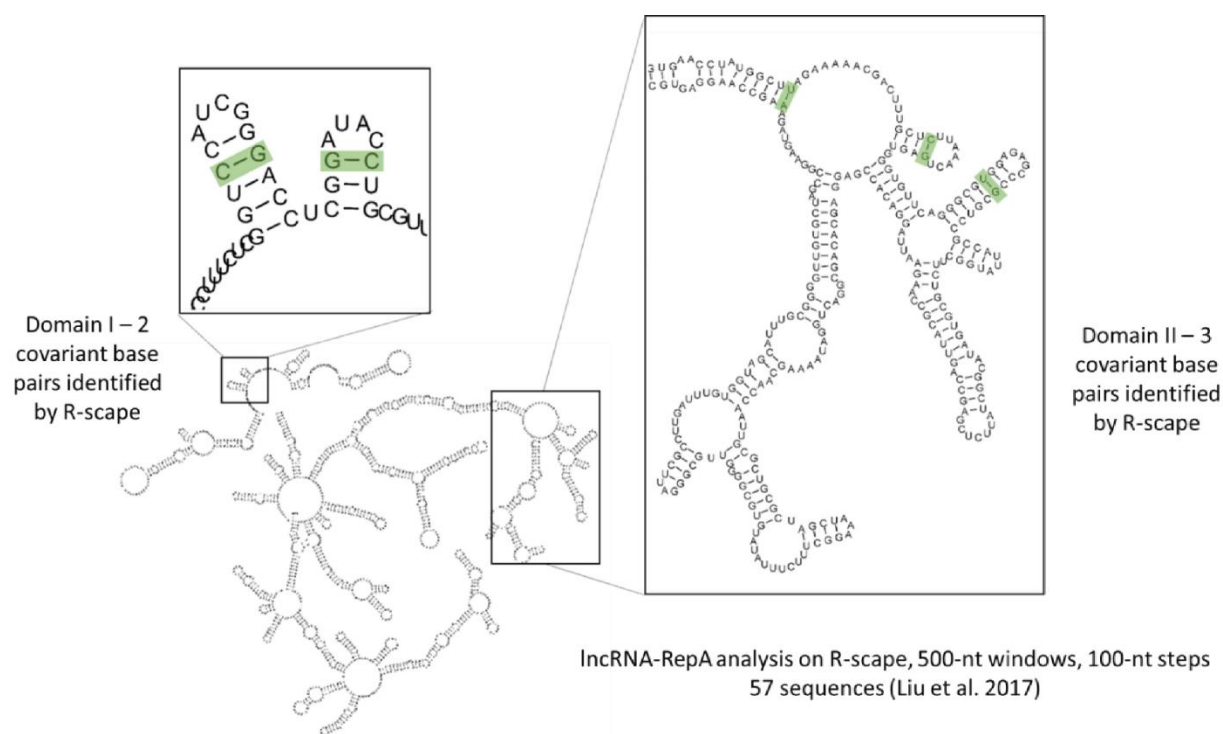
291

## Figure 1

292

293



294

295

296

297

298

299

300

301

302

303

**Figure 2**

304



SAM Riboswitch    SAM Riboswitch + mRNA region

305

306

307

308

309

310

311

**Figure 3**



313

314

315

316

317 **Figure 4**

318



Domain I – 2 covariant base pairs identified by R-scape

Domain II – 3 covariant base pairs identified by R-scape

lncRNA-RepA analysis on R-scape, 500-nt windows, 100-nt steps
57 sequences (Liu et al. 2017)
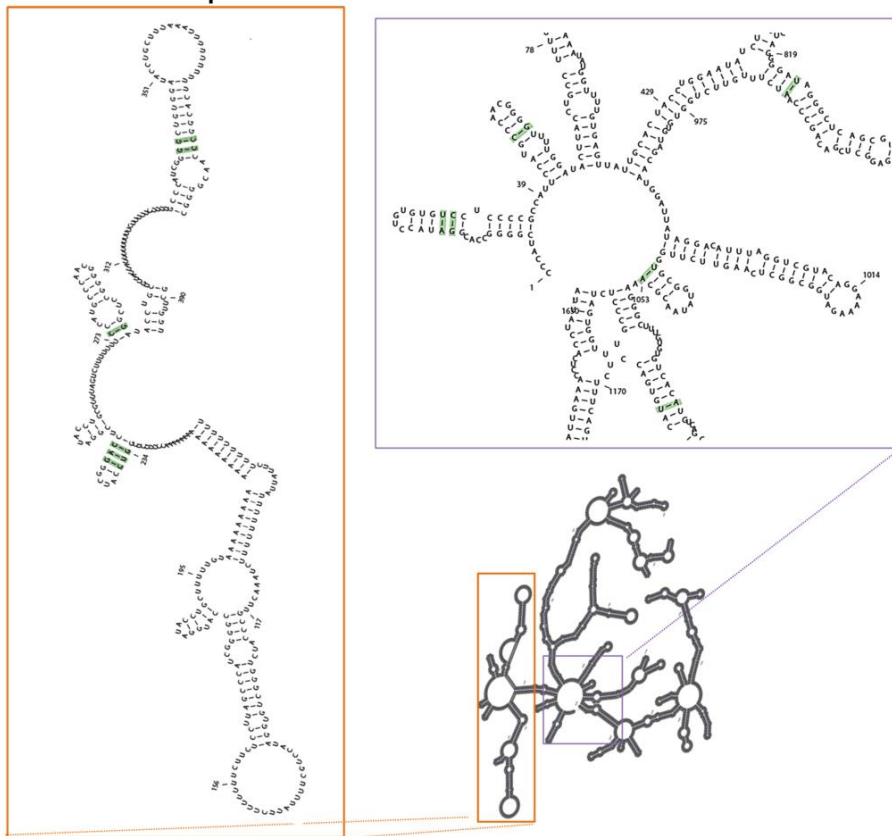
319

320

321

322

323

324

325

326

327

328

329
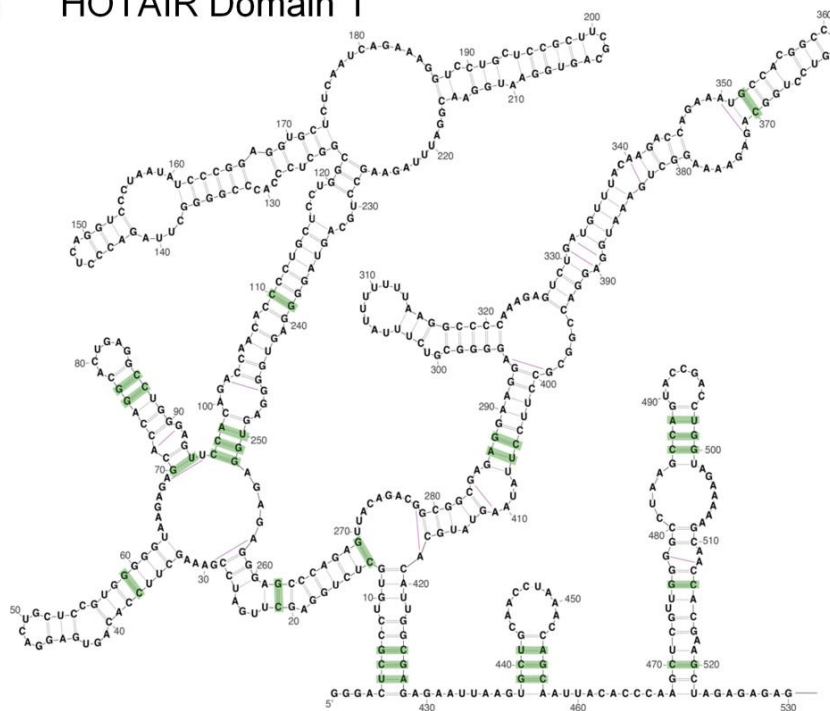
330     **Figure 5**

a     lncRNA-RepA



b     HOTAIR Domain 1



331

## Figure legends

**Figure 1**

Restriction in alignment characteristics (number of sequences, average pairwise sequence identity and phylogenetic diversity) significantly impair R-scape's ability to detect covariation in highly conserved structured RNAs. The percentage of covariant basepairs flagged by R-scape is shown in the graph for each tested RNA alignment.

**Figure 2**

R-scape analysis on the SAM-I riboswitch (RF00162) with and without unstructured mRNA regions in the alignment. (**a**) Sensitivity of R-scape to the presence of adjacent unstructured regions, as a function of the number of sequences in the alignment. (**b**) Influence of an adjacent unstructured region on predicted covariation in the SAM-I riboswitch, using 60 sequences in the alignment. The figure shows the graphical output of each analysis generated by R-scape using R2R drawing notation. Green boxes indicate covariant basepairs. Consensus nucleotide letters are colored according to their sequence conservation in the alignment as given by percent identity thresholds (75% identity in gray; 90% identity in black; 97% identity in red). Individual nucleotides are represented in circles according to their positional conservation in the alignment corresponding to percent occupancy thresholds (50% occupancy in white; 75% occupancy in gray; 90% occupancy in black; 97% occupancy in red).

**Figure 3**

Sliding windows analysis improves R-scape performance on long alignments. In both model cases tested in this study, 7SK RNA (**a**) and Aphthovirus IRES (**b**), R-scape identified four additional base-pairs when the analysis was run in sliding windows.  The consensus secondary

19

354    structure of each RNA is shown in the cartoon form below, and insets above show the

355    covariation predictions for specific domains.  Predicted covariant base pairs are highlighted in

356    green.

357    **Figure 4**

358    R-scape analysis on lncRNA RepA's recently published structure (Liu et al. 2017). The use of an

359    alignment containing 57 sequences was coupled with a sliding windows approach in order to

360    improve covariation analysis on R-scape. The experimentally determined secondary structure of

361    the lncRNA is represented in the figure with insets showing the covariant basepairs (green

362    boxes) identified by R-scape on specific motifs of domain I and domain II of RepA (left and

363    right insets, respectively).
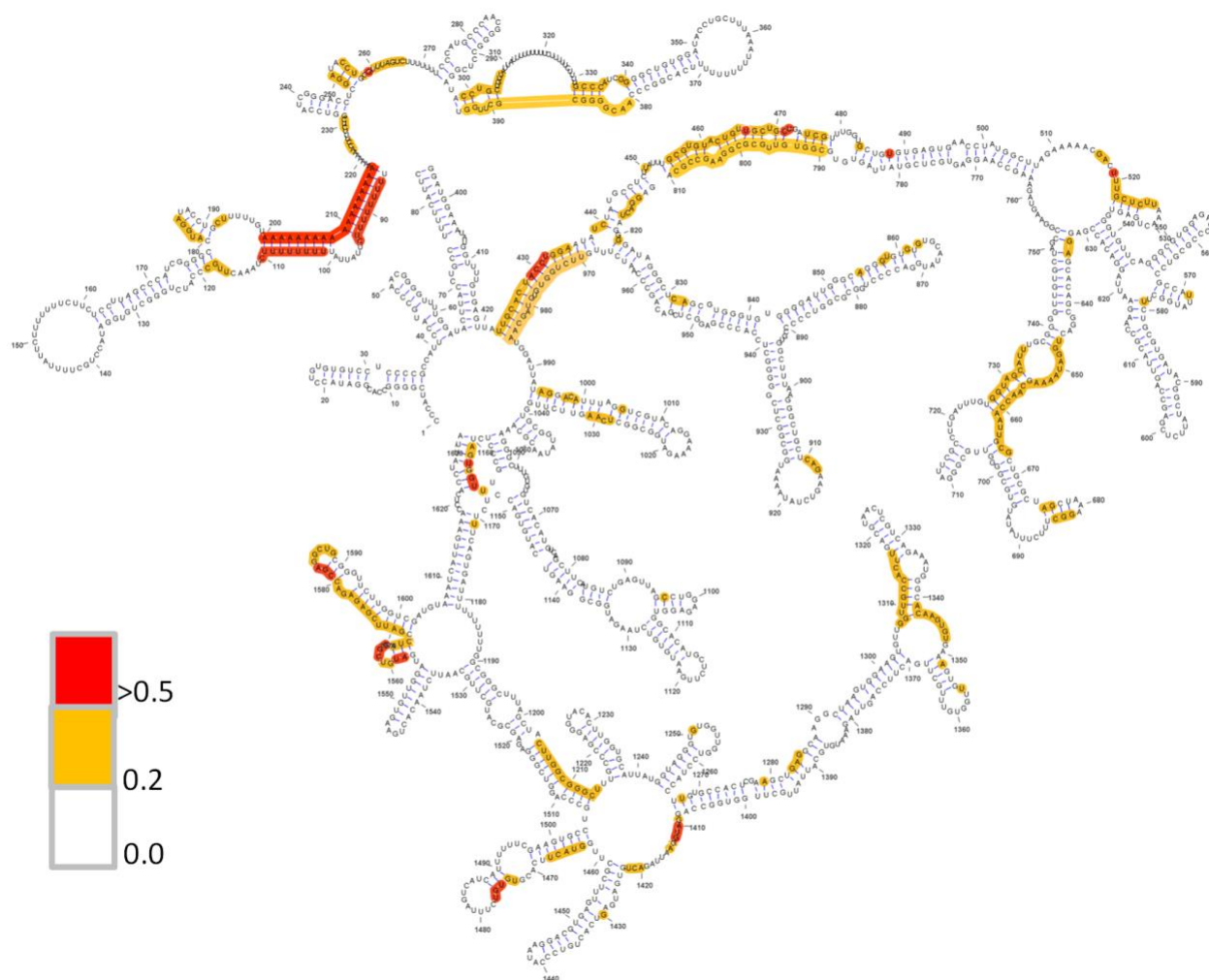
364    **Figure 5**

365    R-scape analysis on lncRNAs RepA and HOTAIR using APC-RAFS as the covariation metric.

366    (**a**) The experimental secondary structure map of full-length lncRNA RepA is shown and

367    covariant basepairs identified on specific motifs by R-scape using APC-RAFS are indicated in

368    green boxes. (**b**) The experimental secondary structure of domain I of HOTAIR is represented in

369    the figure and the covariant basepairs identified by R-scape using APC-RAFS are shown in

370    green boxes.

371

372

373

374

375 **Supplementary Figure 1**

376

| RNA (RFAM ID) | Alignment description | Number of sequences | Average pairwise sequence identity (%) |
|---|---|---|---|
| tRNA (RF00005) | RFAM seed | 954 | 44.44 |
| | Sub-sampled control | 39 | 46.68 |
| | Mammals-only | 39 | 48.26 |
| 5S rRNA (RF 00001) | RFAM seed | 712 | 56.09 |
| | Sub-sampled control | 33 | 72.93 |
| | Mammals-only | 33 | 78.01 |
| 5.8S rRNA (RF 00002) | RFAM seed | 61 | 67.92 |
| | Sub-sampled control | 32 | 76.50 |
| | Mammals-only | 32 | 75.78 |
| Eukarya RNAse P (RF 00009) | RFAM seed | 116 | 49.12 |
| | Sub-sampled control | 46 | 68.50 |
| | Mammals-only | 45 | 67.78 |
| U2 snRNA (RF 00004) | RFAM seed | 208 | 59.35 |
| | Sub-sampled control | 46 | 66.65 |
| | Mammals-only | 46 | 65.76 |
| U5 snRNA (RF 00020) | RFAM seed | 180 | 52.67 |
| | Sub-sampled control | 44 | 81.19 |
| | Mammals-only | 44 | 83.57 |
| Eukarya SS rRNA (RF RF01960) | RFAM seed | 91 | 62.65 |
| | Sub-sampled control | 33 | 64.91 |
| | Mammals-only | 33 | 64.27 |

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

21

394   **Supplementary Figure 2**



395

396

397

398

399

400

401

402 **Supplementary Figure 3**



403

404

405

406

407
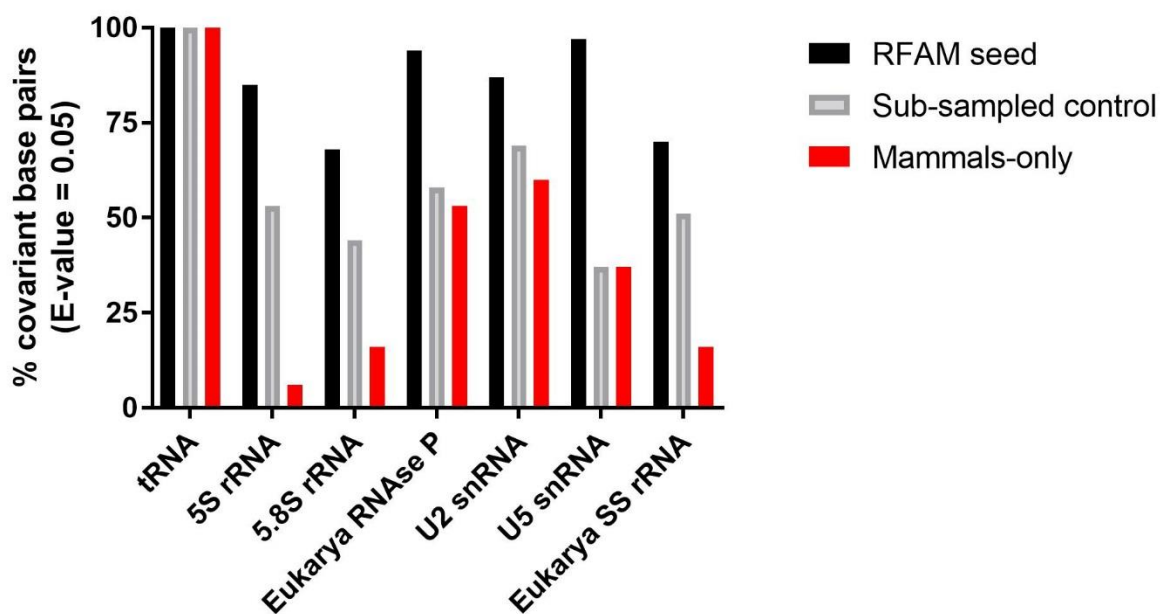
408

409 **Supplementary Figure 4**

410



24

431     **Supplementary Figure 5**

| RNA (RFAM ID) | Alignment description | R-scape search properties (%) Covariation method = APC-GT | | | R-scape search properties (%) Covariation method = APC-RAFS | | |
|---|---|---|---|---|---|---|---|
| | | Sensitivity | PPV | F-measure | Sensitivity | PPV | F-measure |
| tRNA (RF00005) | RFAM seed | 100.0 | 56.76 | 72.41 | 100.00 | 38.89 | 56.00 |
| | Sub-sampled control | 95.24 | 95.24 | 95.24 | 100.00 | 100.00 | 100.00 |
| | Mammals-only | 85.71 | 100.00 | 92.31 | 100.00 | 100.00 | 100.00 |
| 5S rRNA (RF 00001) | RFAM seed | 64.71 | 73.3 | 68.75 | 85.29 | 87.88 | 86.57 |
| | Sub-sampled control | 41.18 | 93.33 | 57.14 | 52.94 | 94.74 | 67.92 |
| | Mammals-only | 0.00 | 0.00 | 0.00 | 5.88 | 100.00 | 11.11 |
| 5.8S rRNA (RF 00002) | RFAM seed | 44.00 | 100.00 | 61.11 | 68.00 | 100.00 | 80.95 |
| | Sub-sampled control | 24.00 | 100.00 | 38.71 | 44.00 | 100.00 | 61.11 |
| | Mammals-only | 0.00 | 0.00 | 0.00 | 16.00 | 100.00 | 27.59 |
| Eukarya RNAse P (RF 00009) | RFAM seed | 79.03 | 100.00 | 88.29 | 93.55 | 81.69 | 87.22 |
| | Sub-sampled control | 12.90 | 88.89 | 22.54 | 58.06 | 100.00 | 73.47 |
| | Mammals-only | 3.64 | 100.0 | 7.02 | 52.73 | 100.00 | 69.05 |
| U2 snRNA (RF 00004) | RFAM seed | 86.67 | 88.64 | 87.64 | 86.67 | 95.12 | 90.70 |
| | Sub-sampled control | 48.89 | 100.00 | 65.67 | 68.89 | 96.88 | 80.52 |
| | Mammals-only | 15.56 | 100.00 | 26.92 | 60.00 | 100.00 | 75.00 |
| U5 snRNA (RF 00020) | RFAM seed | 90.00 | 77.14 | 83.08 | 96.67 | 74.36 | 84.06 |
| | Sub-sampled control | 20.00 | 100.00 | 33.33 | 36.67 | 100.00 | 53.66 |
| | Mammals-only | 3.33 | 100.00 | 6.45 | 36.67 | 100.00 | 53.66 |
| Eukarya SS rRNA (RF RF01960) | RFAM seed | 47.87 | 98.17 | 64.36 | 70.25 | 98.12 | 81.88 |
| | Sub-sampled control | 24.38 | 100.00 | 39.21 | 51.45 | 100.00 | 67.95 |
| | Mammals-only | 1.34 | 100.00 | 2.65 | 15.88 | 98.61 | 27.36 |

432

25

## Supplementary figure legends

### Supplementary Figure 1

Parameters of structural alignments used for the R-scape analysis presented in Figure 1.  RFAM IDs are indicated for each RNA family and the number of sequences and average pairwise sequence identity of each individual alignment (seed alignment, sub-sampled and mammalian sequences) are listed.

### Supplementary Figure 2

Shannon entropy values mapped onto the experimental secondary structure map of lncRNA-RepA (Adapted from Liu et al. 2017). Nucleotides with high Shannon entropy values are represented in red ($> 0.5$) circles; those with medium values (0.2-0.5) are represented in yellow circles. Nucleotides with low Shannon entropy ($< 0.2$) are not highlighted in the map.

### Supplementary Figure 3

Comparison between the APC-GT and APC-RAFS covariation statistics currently implemented in R-scape. (**a**) Difference in Sensitivity (Sen) and Positive Predictive Value (PPV) between APC-GT and APC-RAFS at various E-value thresholds. At the R-scape default E-value, APC-RAFS shows much better sensitivity over APC-GT. (**b**) Same analysis as in (**a**), now varying the number of sequences in the alignments to include a range more commonly found in lncRNA alignments. At a fixed E-value threshold (0.05), APC-RAFS results in superior sensitivity even with 10 sequences in the alignment.

454 **Supplementary Figure 4**

455 R-scape analysis using APC-RAFS as the covariation method on the same alignments used in

456 Figure 1 of the main text. The percentage of covariant basepairs flagged by R-scape as

457 statistically significant is shown in the graph for each tested RNA alignment.

458 **Supplementary Figure 5**

459 R-scape search parameters for the alignments referred to in Figure 1 and Supplementary Figure 1,

460 comparing APC-GT and APC-RAFS covariation statistics. RFAM IDs are indicated for each

461 RNA family. The percent values of sensitivity, positive predictive value and F-measure of each

462 R-scape search were obtained from the analysis output for each alignment (seed alignment, sub-

463 sampled and mammalian sequences) using both covariation methods, APC-GT (R-scape's

464 default) and APC-RAFS.

465

466