

IRIS: A Method for Predicting *in vivo* RNA Secondary Structures Using PARIS Data

Supplementary information

Jiayu Zhou^{1,2}, Pan Li³, Wanwen Zeng^{1,4}, Wenxiu Ma⁵, Zhipeng Lu⁶, Rui Jiang^{1,7}, Qiangfeng Cliff Zhang^{3,*}, Tao Jiang^{8,1,2,*}

¹ *Bioinformatics Division, BNRIST, Tsinghua University, Beijing 100084, China*

² *Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China*

³ *MOE Key Laboratory of Bioinformatics, Beijing Advanced Innovation Center for Structural Biology, Center for Synthetic and Systems Biology, Tsinghua-Peking Center for Life Sciences, School of Life Sciences, Tsinghua University, Beijing 100084, China*

⁴ *College of Software, Nankai University, Tianjin 300071, China*

⁵ *Department of Statistics, University of California, Riverside, CA 92521, USA*

⁶ *Department of Pharmacology and Pharmaceutical Sciences, University of Southern California, CA 90089, USA*

⁷ *Department of Automation, Tsinghua University, Beijing 100084, China*

⁸ *Department of Computer Science and Engineering, University of California, Riverside, CA 92521, USA*

Contents

S1	Important derivation	1
S1.1	Determination of λ in the picking step.....	1
S1.2	Derivation of the log-likelihood of icSHAPE data.....	1
S2	Parameter configuration.....	3
S2.1	Parameters for running STAR	3
S2.2	Parameters for running IRIS.....	3
S3	Additional results	3

S1 Important derivation

S1.1 Determination of λ in the picking step

The parameter λ is used to scale the likelihood distribution $P(\mathbf{Y}|\mathcal{X}_t, \boldsymbol{\alpha}_t)$ to make it comparable with the prior distribution $P(\mathcal{X}_t, \boldsymbol{\alpha}_t)$. Before deriving the parameter λ , we first define the prior and likelihood distribution for a single structure from generated candidate structures as in Equations S1 and S2 according based on Equations 9 and 10:

$$P(\mathbf{X}_c) = \frac{1}{Z} \exp\left(-\frac{1}{\beta} \mathcal{E}(\mathbf{X}_c)\right) \quad (\text{S1})$$

$$P(\mathbf{Y}|\mathbf{X}_c) = \lambda \exp\left(-\lambda(\mathcal{S}_{max} - \mathcal{S}(\mathbf{X}_c))\right) \quad (\text{S2})$$

In this way, the determination of parameter λ is independent of K (the number of structures in the desired ensemble) by solving the equation as in Equation S3:

$$\frac{\max_c P(\mathbf{Y}|\mathbf{X}_c)}{\min_c P(\mathbf{Y}|\mathbf{X}_c)} = \frac{\max_c P(\mathbf{X}_c)}{\min_c P(\mathbf{X}_c)} \quad (\text{S3})$$

Then, we can easily conclude that the maximum value of $P(\mathbf{Y}|\mathbf{X}_c)$ is obtained when $\mathcal{S}(\mathbf{X}_c) = \mathcal{S}_{max}$ and the minimum value is obtained when $\mathcal{S}(\mathbf{X}_c) = \mathcal{S}_{min}$, where \mathcal{S}_{max} and \mathcal{S}_{min} respectively denote the maximum and minimum PARIS support among the candidate structures. Finally, through logarithm transform and simple mathematical manipulation, we can get the parameter λ as in Equation S4:

$$\lambda = \frac{\max_c \{\log P(\mathbf{X}_c)\} - \min_c \{\log P(\mathbf{X}_c)\}}{\mathcal{S}_{max} - \mathcal{S}_{min}} \quad (\text{S4})$$

S1.2 Derivation of the log-likelihood of icSHAPE data

First, the scores of icSHAPE data for paired and unpaired bases are fitted separately by using the Beta distribution implemented in SciPy as follows:

$$\begin{aligned} \mathcal{B}_{\text{paired}}(z_i) &= \text{Beta}(z_i | \alpha = 0.158, \beta = 0.846) \\ \mathcal{B}_{\text{unpaired}}(z_i) &= \text{Beta}(z_i | \alpha = 0.209, \beta = 0.365) \end{aligned} \quad (\text{S5})$$

Then, based on the assumption that the icSHAPE scores of each base on the RNA are independent of each other, the log-likelihood of icSHAPE scores \mathbf{z} of the RNA observed from the predicted set $\tilde{\mathcal{X}}, \tilde{\boldsymbol{\alpha}}$ can be decomposed by Equation S6:

$$\log P(\mathbf{z}|\tilde{\mathcal{X}}, \tilde{\boldsymbol{\alpha}}) = \sum_{i=1}^n \log P(z_i|\tilde{\mathcal{X}}, \tilde{\boldsymbol{\alpha}}) \quad (\text{S6})$$

Next, the probability of observing z_i given the ensemble is a mixture of observing z_i from each structure, *i.e.*,

$$\log P(z_i|\tilde{\mathcal{X}}, \tilde{\boldsymbol{\alpha}}) = \log \left(\sum_{k=1}^K \tilde{\alpha}_k P(z_i|\tilde{\mathbf{X}}_k) \right) \quad (\text{S7})$$

The probability $P(z_i|\tilde{\mathbf{X}}_k)$ is only dependent on whether base i is paired in structure k , that is

$$P(z_i|\tilde{\mathbf{X}}_k) = \left(\sum_{j=1}^n \tilde{x}_{kij} \right) \mathcal{B}_{\text{paired}}(z_i) + \left(1 - \left(\sum_{j=1}^n \tilde{x}_{kij} \right) \right) \mathcal{B}_{\text{unpaired}}(z_i) \quad (\text{S8})$$

Note that $\sum_{j=1}^n \tilde{x}_{kij}$ must be 0 or 1 due to the fact that each base can either form a base pair or be unpaired. Here, we expand the base-pairing probability of base i defined in the main content as in Equation S9:

$$b_i = \sum_{j=1}^n b_{ij} = \sum_{j=1}^n \sum_{k=1}^K \tilde{\alpha}_k \tilde{x}_{kij} = \sum_{k=1}^K \tilde{\alpha}_k \sum_{j=1}^n \tilde{x}_{kij} \quad (\text{S9})$$

Finally, we can easily derive Equation 2 by substituting equations S7, S8 and S9 into Equation S6.

S2 Parameter configuration

S2.1 Parameters for running STAR

As mentioned in the main text, because STAR was originally designed for analyzing alternative splicing in RNA-Seq data, we modify the default parameters of STAR to remove the biases that take into account certain natures of slicing as **Table S1**.

Table S1 Special parameters used in STAR for mapping PARIS reads

Parameter	Value	Parameter	Value
outSJfilterReads	Unique	outSJfilterOverhangMin	5 5 5 5
alignIntronMin	10	outSJfilterCountUniqueMin	1 1 1 1
scoreGapNoncan	0	outSJfilterCountTotalMin	1 1 1 1
scoreGapGCAG	0	outSJfilterDistToOtherSJmin	0 0 0 0
scoreGapATAC	0		

S2.2 Parameters for running IRIS

Three parameters need to be configured before running IRIS, which are the range of length k for short stems, the fraction of PARIS support in Y as the threshold for filtering stems, and the number of clusters C . Empirically, the number of clusters is set to 100 for all 11 RNAs in our benchmarking experiments. But the other two parameters vary according to the length of the input RNA. For RNAs less than 160 nt in length, k takes values from 3 to 6 and the fraction threshold is 0.5. For RNAs longer than 160 nt but shorter than 300 nt, k is from 4 to 7 and the fraction threshold is 0.75. Finally, for RNAs longer than 300 nt, the value of k is 4 to 7 and the fraction threshold is 0.90.

S3 Additional results

We show the results concerning RMRP as an additional case, as shown in **Figure S1**. From Figure S1A, we note that all structures predicted by IRIS contain a long continuous stem structure with the center located near 150 nt. The reason is that this

region has PARIS support enriched (Figure S1C), and it can be verified by evidence of evolutionary conservation (marked by the black arrow in Figure S1B).

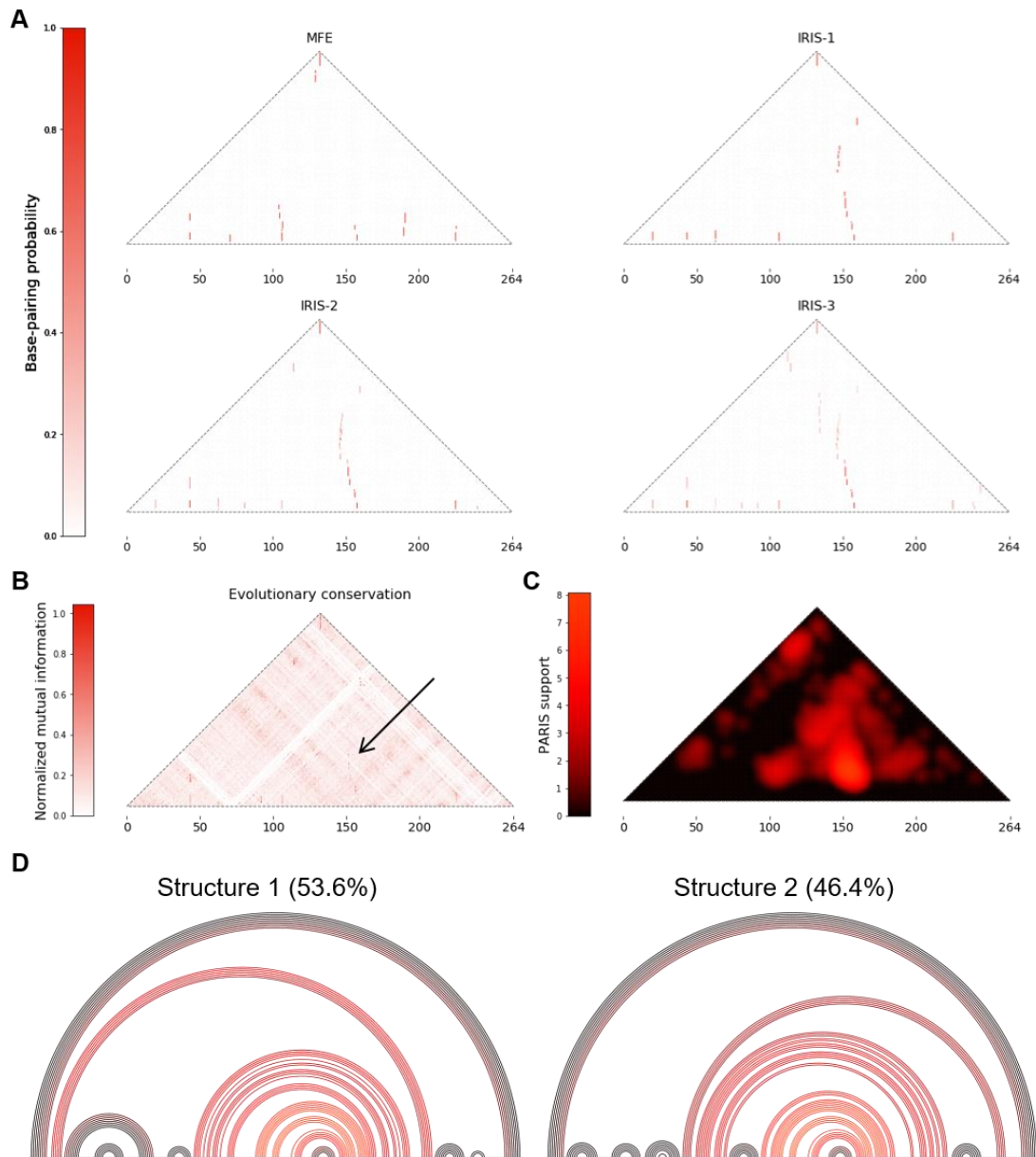


Figure S1 Results concerning the RMRP

A. The base-pairing probability matrices of the predicted ensemble from MFE, IRIS-1, IRIS-2 and IRIS-3. **B.** The matrix of normalized mutual information. **C.** The matrix of PARIS support. **D.** The ensemble of two representative structures predicted by IRIS-2.