

Results

Basic Notation and Definitions

Consider a non-coding RNA molecule with N nucleotides. Counting from 5' to 3', we denote the primary structure by

$$p = (p_1, p_2, \dots, p_N), \text{ where } p_i \in \{A, G, C, U\}, i = 1, \dots, N \quad (1)$$

and the secondary structure by

$$s = \{(j, k) : \text{nucleotides } j \text{ and } k \text{ are paired}, 1 \leq j < k \leq N\} \quad (2)$$

and we define the *segment* at *location* i to be

$$P_i = (P_{i,1}, P_{i,2}, P_{i,3}, P_{i,4}) = (p_i, p_{i+1}, p_{i+2}, p_{i+3}) \quad (3)$$

There is no particular reason for using segments of length four, and in fact all qualitative conclusions are identical with segment lengths three, four, or five, and, for that matter, whether or not we include the G·U wobble pairs. We chose to include them.

Which segments are viable complementary pairs to P_i ? The only constraint on location is that an RNA molecule cannot form a loop of two or fewer nucleotides. Let A_i be the set of all segments that are potential pairs of P_i :

$$A_i = \{P_j : 1 \leq j \leq i - 7 \text{ or } i + 7 \leq j \leq N - 3\} \quad (4)$$

We can now define the *local ambiguity function*,

$$a(p) = (a_1(p), \dots, a_{N-3}(p))$$

which is a vector-valued function of the primary structure p . The vector has one component, $a_i(p)$, for each segment P_i , which is, simply, the number of feasible segments that are complementary to P_i :

$$\begin{aligned} a_i(p) &= \#\{P \in A_i : P \text{ and } P_i \text{ are complementary}\} \\ &= \#\{P_j \in A_i : (P_{i,k}, P_{j,5-k}) \in \{(A, U), (U, A), (G, C), (C, G), (G, U), (U, G)\}, \\ &\quad k = 1, \dots, 4\} \end{aligned} \quad (5)$$

We want to explore the relationship between ambiguity and secondary structure. We can do this conveniently, on a molecule-by-molecule basis, by introducing another vector-valued function, this time depending only on a purported secondary structure. Specifically, the new function assigns a descriptive label to each location (i.e. each nucleotide), determined by whether the segment at the given location is fully paired, partially paired, or fully unpaired.

Formally, given a secondary structure s , as defined in Eq (2), and a location $i \in \{1, 2, \dots, N - 3\}$, let $f_i(s)$ be the number of nucleotides in P_i that are paired under s :

$$f_i(s) = \#\{j \in P_i : (j, k) \in s \text{ or } (k, j) \in s, \text{ for some } 1 \leq k \leq N\} \quad (6)$$

Evidently, $0 \leq f_i(s) \leq 4$. The “paired nucleotides function” is then the vector-valued function of secondary structure defined as $f(s) = (f_1(s), \dots, f_{N-3}(s))$. Finally, we use f to distinguish three types of locations (and hence three types of segments): location i will be labeled

$$\begin{cases} \text{single} & \text{if } f_i(s) = 0 \\ \text{double} & \text{if } f_i(s) = 4 \\ \text{transitional} & \text{if } 0 < f_i(s) < 4 \end{cases} \quad i = 1, 2, \dots, N - 3 \quad (7)$$

A First Look at the Data: Shuffling Nucleotides

Our goals are to explore connections between ambiguities and basic characteristics of RNA families, as well as the changes in these relationships, if any, when using comparative, as opposed to MFE, secondary structures. For each molecule and each location i , the segment at i has been assigned a “local ambiguity” $a_i(p)$ that depends only on the primary structure, and a label (*single*, *double*, or *transitional*) that depends only on the secondary structure. Since the local ambiguity, by itself, is strongly dependent on the length of the molecule, and possibly on other intrinsic properties, we propose two *relative* ambiguity indexes: T-S and D-S, each of which depends on both the primary (p) and purported secondary (s) structures. Formally,

$$d_{\text{T-S}}(p, s) = \frac{\sum_{j=0}^{N-3} a_j(p) c_j^{\text{tran}}(s)}{\sum_{j=0}^{N-3} c_j^{\text{tran}}(s)} - \frac{\sum_{j=0}^{N-3} a_j(p) c_j^{\text{single}}(s)}{\sum_{j=0}^{N-3} c_j^{\text{single}}(s)} \quad (8)$$

where we have used c_i^{tran} and c_i^{single} for indicating whether location i is *transitional* or *single* respectively. In other words, for each $i = 1, 2, \dots, N - 3$

$$c_i^{\text{tran}}(s) = \begin{cases} 1, & \text{if location } i \text{ is } \textit{transitional} \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

$$c_i^{\text{single}}(s) = \begin{cases} 1, & \text{if location } i \text{ is } \textit{single} \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

In short, the T-S ambiguity index is the difference in the averages of the local ambiguities at *transitional* sites and *single* sites. Since, as noted earlier, the local ambiguity at every *double* location, i , is likely to be at least one ($a_i \geq 1$), the D-S index involves only those *double* and *single* locations, j , for which $a_j \geq 1$:

$$d_{\text{D-S}}(p, s) = \frac{\sum_{j=0}^{N-3} a_j(p) \hat{c}_j^{\text{double}}(p, s)}{\sum_{j=0}^{N-3} \hat{c}_j^{\text{double}}(p, s)} - \frac{\sum_{j=0}^{N-3} a_j(p) \hat{c}_j^{\text{single}}(p, s)}{\sum_{j=0}^{N-3} \hat{c}_j^{\text{single}}(p, s)} \quad (11)$$

Here, $\hat{c}_i^{\text{double}}$ and $\hat{c}_i^{\text{single}}$ indicates those *double* and *single* locations with at least one pair: for each $i = 1, 2, \dots, N - 3$,

$$\hat{c}_i^{\text{double}}(p, s) = \begin{cases} 1, & \text{if location } i \text{ is } \textit{double} \text{ and } a_i(p) \geq 1 \\ 0, & \text{otherwise} \end{cases} \quad (12)$$

$$\hat{c}_i^{\text{single}}(p, s) = \begin{cases} 1, & \text{if location } i \text{ is } \textit{single} \text{ and } a_i(p) \geq 1 \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

Thinking kinetically, we might expect to find relatively small values of the T-S and D-S indexes ($d_{\text{T-S}}$ and $d_{\text{D-S}}$) in RNA families which fold and operate as individual entities—what we call here the unbound RNA. One reason for believing this is that larger numbers of partial matches for a given sequence in or around a stem would likely interfere with the *nucleation* of the native stem structure, potentially disrupting folding, and nucleation appears to be a critical and perhaps even rate-limiting step in folding. Indeed, the experimental literature [4–7] has long suggested that stem formation in RNA molecules is a two-step process. When forming a stem, there is usually a slow nucleation step, resulting in a few consecutive base pairs at a nucleation point, followed by a fast zipping step. Since the ambiguity indexes are functions of the secondary structure, s , our expectations are made, implicitly, under the assumption that s is

correct. For the time being we will focus only on comparative structures, returning later to the questions about MFE structures raised in the *Introduction*.

How are we to gauge d_{T-S} and d_{D-S} , and how are we to compare their values across different RNA families? Consider the following experiment: for a given RNA molecule we create a “surrogate” which has the same nucleotides, and in fact the same counts of *all* four-tuple segments as the original molecule, but is otherwise ordered randomly. If ACCU appeared eight times in the original molecule, then it appears eight times in the surrogate, and the same can be said of all sequences of four successive nucleotides—the frequency of each of the 4^4 possible segments is preserved in the surrogate. If we also preserve the locations of the *transitional*, *double*, and *single* labels, then we can compute new values for d_{T-S} and d_{D-S} , say \tilde{d}_{T-S} and \tilde{d}_{D-S} , from the surrogate. If we produce many surrogate sequences then we will get a sampling of surrogate values, \tilde{d}_{T-S} and \tilde{d}_{D-S} , to which we can compare d_{T-S} and d_{D-S} , respectively. We made several experiments of this type—for each of the two indexes, T-S and D-S, within each of the seven RNA families (Group I and Group II Introns, tmRNA, SRP RNA, RNase P, and 16s and 23s rRNA).

To make this precise, consider an RNA molecule with primary structure p and comparative secondary structure s . Construct a segment “histogram function,” $\mathcal{H}(p)$, which is the vector of the number of times that each of the 4^4 possible segments appears in p .¹ Let $\mathcal{P}(p)$ be the set of all permutations of the ordering of nucleotides in p , and let $\mathcal{E}(p) \subseteq \mathcal{P}(p)$ be the subset of permutations that preserve the frequencies of four-tuples:

$$\mathcal{E}(p) = \{p' \in \mathcal{P}(p) : \mathcal{H}(p') = \mathcal{H}(p)\}$$

Clever algorithms exist for efficiently drawing independent samples from the uniform distribution on \mathcal{E} —see [1–3]. Let $p^{(1)}, \dots, p^{(K)}$ be K such samples, and let $d_{T-S}(p^{(1)}, s), \dots, d_{T-S}(p^{(K)}, s)$ and $d_{D-S}(p^{(1)}, s), \dots, d_{D-S}(p^{(K)}, s)$ be the corresponding T-S and D-S ambiguity indexes. The results are not very sensitive to K , nor to the particular sample, provided that K is large enough. We used $K = 10,000$. Finally, let $\alpha_{T-S}(p, s) \in [0, 1]$ and $\alpha_{D-S}(p, s) \in [0, 1]$ be the left-tail empirical probabilities, under the distributions defined by the ensembles

$$d_{T-S}(p, s), d_{T-S}(p^{(1)}, s), \dots, d_{T-S}(p^{(K)}, s) \quad (14)$$

$$\text{and} \quad d_{D-S}(p, s), d_{D-S}(p^{(1)}, s), \dots, d_{D-S}(p^{(K)}, s) \quad (15)$$

of choosing an ambiguity index less than or equal to $d_{T-S}(p, s)$ and $d_{D-S}(p, s)$, respectively:

$$\alpha_{T-S}(p, s) = \frac{1 + \#\{k \in \{1, \dots, K\} : d_{T-S}(p^{(k)}, s) \leq d_{T-S}(p, s)\}}{1 + K} \quad (16)$$

$$\alpha_{D-S}(p, s) = \frac{1 + \#\{k \in \{1, \dots, K\} : d_{D-S}(p^{(k)}, s) \leq d_{D-S}(p, s)\}}{1 + K} \quad (17)$$

In essence, each α score is a self-calibrated ambiguity index.

It is tempting to interpret $\alpha_{T-S}(p, s)$ as a p-value from a conditional hypothesis test: Given s and \mathcal{H} , test the null hypothesis that $d_{T-S}(p, s)$ is statistically indistinguishable from $d_{T-S}(p', s)$, where p' is a random sample from \mathcal{E} . If the alternative hypothesis were that $d_{T-S}(p, s)$ is too small to be consistent with the null, then the null is rejected in favor of the alternative with probability $\alpha_{T-S}(p, s)$. The problem with this interpretation, and the analogous one for $\alpha_{D-S}(p, s)$, is that this null hypothesis violates the observation that given \mathcal{H} there is information in s about p , whereas $p^{(1)}, \dots, p^{(K)}$ are independent of s given \mathcal{H} . In other words, $d_{T-S}(p, s)$ and $d_{T-S}(p', s)$ have different conditional distributions given s and \mathcal{H} , in direct contradiction to the null hypothesis.

¹ \mathcal{H} is made a vector by fixing an arbitrary ordering of all possible segments.

A larger problem is that there is no reason to believe the alternative; we are more interested in *relative* than absolute ambiguity indexes. Thinking of $\alpha_{\text{T-S}}(p, s)$ and $\alpha_{\text{D-S}}(p, s)$ as calibrated intra-molecular indexes, we want to know how they vary across RNA families, and whether these variations depend on the differences between comparative and MFE structures.

Nevertheless, $\alpha_{\text{T-S}}(p, s)$ and $\alpha_{\text{D-S}}(p, s)$ are useful statistics for exploratory analysis. Table 1 provides summary data about the α scores for each of the seven RNA families. For each molecule in each family we use the primary structure and the comparative secondary structure, and $K=10,000$ samples from \mathcal{E} , to compute individual T-S and D-S α scores (Eqs 16 and 17). Keeping in mind that a smaller value of α represents a smaller *calibrated* value of the corresponding ambiguity index, $d(p, s)$, there is evidently a disparity between ambiguity indexes of RNA molecules that form ribonucleoproteins and those of RNA that are already active as individual molecules. As a group, the unbound molecules have systematically lower ambiguity indexes. As already noted, this observation is consistent with, and in fact anticipated by, the kinetic point of view. Shortly, we will support this observation with ROC curves and rigorous hypothesis tests.

family	number molecules	median length	median $\alpha_{\text{T-S}}$	median $\alpha_{\text{D-S}}$
Group I Introns	116	451	0.432	0.908
Group II Introns	34	990	0.181	0.761
tmRNA	404	363	0.926	0.988
SRP RNA	346	274	0.790	0.967
RNase P	407	330	0.925	0.985
16s rRNA	279	1512	0.938	1.000
23s rRNA	59	2913	1.000	1.000

Table 1. Comparative Secondary Structures: calibrated ambiguity indexes, by RNA family. The number of molecules, the median length (number of nucleotides), and the median α scores for T-S and D-S ambiguity indexes (Eqs 16 and 17) for each of the seven RNA families studied. RNA molecules from the first two families are active as single molecules (unbound); the remaining five are bound in ribonucleoproteins. Unbound RNA molecules have lower ambiguity indexes.

Does the MFE structure similarly separate single-entity RNA molecules from those that form ribonucleoproteins? A convenient way to explore this question is to recalculate and recalibrate the ambiguity indexes of each molecule in each of the seven families, but using the MFE in place of the comparative secondary structures. The results are summarized in Table 2. By comparison to the results shown from Table 1, the separation of unbound from bound molecules nearly disappears under the MFE secondary structures. Possibly, the comparative structures, as opposed to the MFE structures, better anticipate the need to avoid energy traps in the folding landscape. Here too we will soon revisit the data using ROC curves and proper hypothesis tests.

Formal Statistical Analyses

The T-S and D-S ambiguity indexes ($d_{\text{T-S}}(p, s)$ and $d_{\text{D-S}}(p, s)$) are intra-molecular measures of the extent to which segments that lie in and around sub-sequences that are double-stranded have more Watson-Crick and wobble pairings than segments within single-stranded regions. As such, they depend on both p and any purported secondary structure, s . Based on calibrated versions of these indexes ($\alpha_{\text{T-S}}(p, s)$ and $\alpha_{\text{D-S}}(p, s)$) and employing the comparative secondary structure for s we found support for the idea that non-coding RNA molecules which are active as individual entities are more likely

family	number molecules	median length	median α_{T-S}	median α_{D-S}
Group I Introns	116	451	0.833	0.994
Group II Introns	34	990	0.841	0.997
tmRNA	404	363	0.867	0.984
SRP RNA	346	274	0.803	0.953
RNase P	407	330	0.955	0.998
16s rRNA	279	1512	0.982	1.000
23s rRNA	59	2913	1.000	1.000

Table 2. MFE Secondary Structures: calibrated ambiguity indexes, by RNA family.

Identical to Table 1, except that the ambiguity indexes and their calibrations are calculated using the MFE secondary structures rather than comparative analyses. There is little evidence in the MFE secondary structures for lower ambiguity indexes among unbound RNA molecules.

to have small ambiguity indexes than RNA molecules destined to operate as part of ribonucleoproteins. Furthermore, the difference appears to be sensitive to the approach used for identifying secondary structure—there is little, if any, evidence in the MFE secondary structures for lower ambiguities among unbound molecules.

These qualitative observations can be used to formulate precise statistical hypothesis tests. Many tests come to mind, but perhaps the simplest and most transparent are based on nothing more than the molecule-by-molecule signs of the ambiguity indexes. Whereas ignoring the actual values of the indexes is inefficient in terms of information, and probably also in the strict statistical sense, tests based on signs require very few assumptions and are, therefore, more robust to model misspecification. All of the p-values that we will report are based on the hypergeometric distribution, which arises as follows.

We are given a population of M molecules, $m = 1, \dots, M$, each with a binary outcome measure $B_m \in \{-1, +1\}$. There are two subpopulations of interest: the first M_1 molecules make up population 1 and the next M_2 molecules make up population 2; $M_1 + M_2 = M$. We observe n_1 plus values in population 1 and n_2 in population 2

$$n_1 = \#\{m \in \{1, 2, \dots, M_1\} : B_m = +1\} \quad (18)$$

$$n_2 = \#\{m \in \{M_1 + 1, M_1 + 2, \dots, M\} : B_m = +1\} \quad (19)$$

We suspect that population 1 has less than its share of plus ones, meaning that the $n_1 + n_2$ population of plus ones was not randomly distributed among the M molecules. To be precise, let N be the number of plus ones that appear from a draw, without replacement, of M_1 samples from B_1, \dots, B_M . Under the null hypothesis, H_o , n_1 is a sample from the hypergeometric distribution on N :

$$\mathbb{P}\{N = n\} = \frac{\binom{M_1}{n} \binom{M_2}{n_1+n_2-n}}{\binom{M}{n_1+n_2}} \quad \max\{0, n_1 + n_2 - M_2\} \leq n \leq \min\{n_1 + n_2, M_1\} \quad (20)$$

The alternative hypothesis, H_a , is that n_1 is too small to be consistent with H_o , leading to a left-tail test with p-value $\mathbb{P}\{N \leq n_1\}$ (which can be computed directly or using a statistical package, e.g. *hypergeom.cdf* in *scipy.stats*).

It is by now well recognized that p-values should never be the end of the story. One reason is that *any* departure from the null hypothesis in the direction of the alternative, no matter how small, is doomed to be statistically significant, with arbitrarily small p-value, once the sample size is sufficiently large. In addition to reporting p-values, we

will also display estimated ROC curves, summarizing performance of a related classification problem. There are eight such problems: Under each of the comparative and MFE secondary structures, how well can we use $d_{T-S}(p, s)$ or $d_{D-S}(p, s)$ to classify a given molecule as bound or unbound? And, given $d_{T-S}(p, s)$ or $d_{D-S}(p, s)$ for a molecule that is known to be bound or unbound, how well can we classify the secondary structure, s , as coming from a comparative analysis versus an MFE analysis?

Single-entity RNA Molecules versus Protein-RNA Complexes

Consider an RNA molecule, m , selected from one of the seven families in our data set, with primary structure p and secondary structure s , computed by comparative analysis. Given only the T-S ambiguity index of m (i.e. given only $d_{T-S}(p, s)$), how accurately could we classify m as being unbound (i.e. from the Group I or Group II Introns) as opposed to bound (i.e. from one of the five families tmRNA, SRP RNA, RNase P, 16s rRNA or 23s rRNA)? The foregoing exploratory analysis suggests constructing a classifier that declares a molecule to be ‘unbound’ when $d_{T-S}(p, s)$ is small, e.g. $d_{T-S}(p, s) < t$, where the threshold t governs the familiar trade off between rates of “true positives” (an unbound m is declared ‘unbound’) and “false positives” (a bound m is declared ‘unbound’). Small values of t favor low rates of false positives at the price of low rates of true positives, whereas large values of t favor high rates of true-positives at the price of high rates of false positives. Since for each molecule m we have both the correct classification (unbound or bound) and the statistic d , we can estimate the ROC performance of our threshold classifier by plotting the empirical values of the pair

$$(\# \text{ false positives}, \# \text{ true positives})$$

for each value of t . The ROC curve for the two-category (unbound versus bound) classifier based on thresholding $d_{T-S}(p, s) < t$ is shown in the upper-left panel of Figure 1. Also shown is the estimated area under the curve (AUC=0.81), which has a convenient and intuitive interpretation, as it is equal to the probability that for two randomly selected molecules, m from the unbound population and m' from the bound population, the T-S ambiguity index of m will be smaller than the T-S ambiguity index of m' .

As mentioned earlier, we can also associate a traditional p-value to the problem of separating unbound from bound molecules, based on the T-S ambiguity indexes. We consider only the signs (positive or negative) of these indexes, and then test whether there are fewer than expected positive indexes among the unbound, as opposed to the bound, population. This amounts to computing $\mathbb{P}\{N \leq n_1\}$ from the hypergeometric distribution—Eq (20). The relevant statistics can be found in Table 3, under the column labels **#mol’s** and **# $d_{T-S} > 0$** . Specifically, $M_1 = 116 + 34 = 150$ (number of unbound molecules), $M_2 = 404 + 346 + 407 + 279 + 59 = 1495$ (number of bound molecules), $n_1 = 50 + 8 = 58$ (number of positive T-S indexes among unbound molecules) and $n_2 = 368 + 269 + 379 + 210 + 53 = 1279$. The resulting p-value, $1.2 \cdot 10^{-34}$, is essentially zero, meaning that the positive T-S indexes are not distributed proportional to the sizes of the unbound and bound populations, which is by now obvious in any case. To repeat our caution, small p-values conflate sample size with effect size, and for that reason we have chosen additional ways, using permutations as well as classifications, to look at the data.

The comparative secondary structure of an RNA molecule, when combined with its primary structure, can be used to construct a measure—the T-S ambiguity index—which distinguishes unbound from bound molecules with good accuracy. Can the same can be said for the D-S index? Yes, albeit with lower accuracy. To demonstrate, we followed the identical procedure, except that we assigned the index

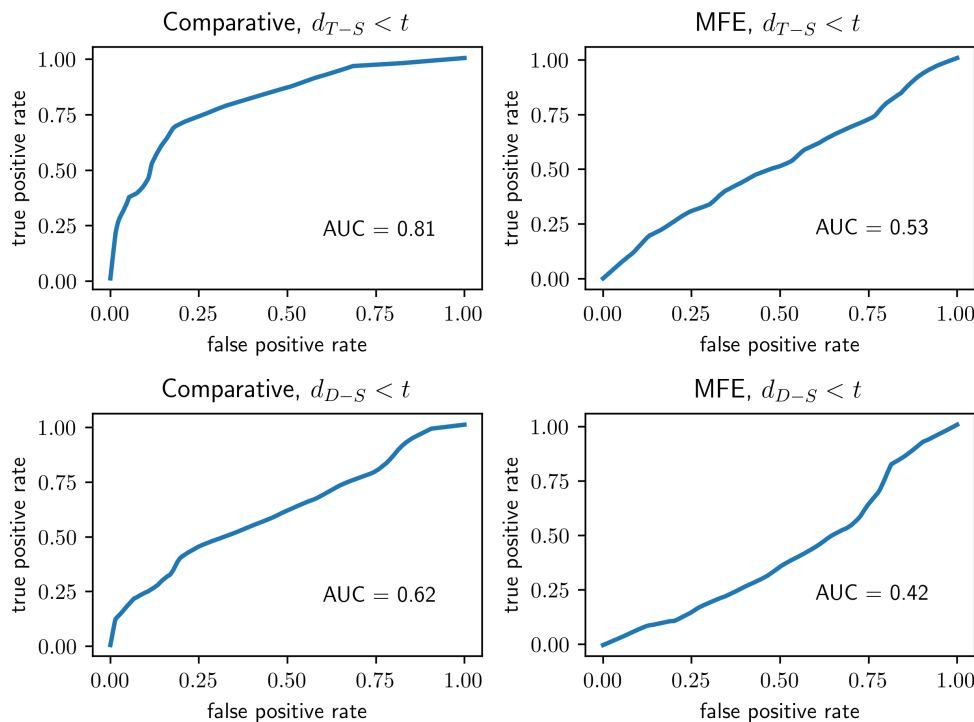


Fig 1. Bound or Unbound? ROC performance of classifiers based on thresholding T-S and D-S ambiguity indexes. Small values are taken as evidence for molecules that are active as single entities (unbound), as opposed to parts of ribonucleoproteins (bound). Classifiers in the left two panels use comparative secondary structures to compute ambiguity indexes; those on the right use (approximate) minimum free energies. In each of the four experiments, a conditional p-value was also calculated, based only on the signs of the indexes and the null hypothesis that positive indexes are distributed randomly among molecules of all types as opposed to the alternative that positive indexes are more typically found among families of bound RNA. Under the null hypothesis, the test statistic is hypergeometric—see Eq 20. *Upper Left:* $p = 1.2 \times 10^{-34}$; *Lower Left:* $p = 7.3 \times 10^{-8}$; *Upper Right:* $p = 0.02$; *Lower Right:* $p = 0.92$. In considering these extreme p-values, it is perhaps worth re-emphasizing the points made about the interpretation of p-values in the paragraph following Eq 20. (These ROC curves and those in Figure 2 were lightly smoothed by the method known as “Locally Weighted Scatterplot Smoothing,” e.g. with the python command `Y=lowess(Y, X, 0.1, return_sorted=False)` coming from `statsmodels.nonparametric.smoothers_lowess`.)

$d_{D-S}(p, s)$ rather than $d_{T-S}(p, s)$ to each molecule. The ROC curve (with area 0.65) is shown in the lower-left panel of Figure 1. The hypergeometric test, based on the signs of the D-S indexes (positive counts, for each RNA family, can be found in Table 3) has p-value $7.3 \cdot 10^{-8}$. Finally, the right-hand panels in Figure 1 mirror the left-hand panels, except that all secondary structures were computed by (approximately) minimizing free energy rather than comparative analysis. The classification results are substantially less convincing and the p-values substantially higher.

family	#mol's	# $d_{T-S} > 0$	# $d_{D-S} > 0$	# $d_{\tilde{T}-\tilde{S}} > 0$	# $d_{\tilde{D}-\tilde{S}} > 0$
Group I Introns	116	50	100	94	114
Group II Introns	34	8	21	27	33
tmRNA	404	368	396	358	396
SRP RNA	346	269	314	264	302
RNase P	407	379	393	377	404
16s rRNA	279	210	251	254	278
23s rRNA	59	53	56	54	58

Table 3. Numbers of Positive Ambiguity Indexes, by family. #mol's: number of molecules; # $d_{T-S} > 0$ and # $d_{D-S} > 0$: numbers of positive T-S and D-S ambiguity indexes, secondary structures computed by *comparative analysis*; # $d_{\tilde{T}-\tilde{S}} > 0$ and # $d_{\tilde{D}-\tilde{S}} > 0$: numbers of positive T-S and D-S ambiguity indexes, secondary structures computed by *minimum free energy*.

Comparative Analysis versus Minimum Free Energy

As we have just seen, ambiguity indexes based on MFE secondary structures, as opposed to comparative secondary structures, do not make the same stark distinction between single and bound RNA molecules. To explore this a little further, we can turn the analyses of the previous paragraphs around and ask to what extent knowledge of the group of a molecule (unbound or bound) and its ambiguity index (e.g. the T-S ambiguity index) is sufficient to predict the source of the secondary structure—comparative or free energy?

Interestingly, there is a sharp difference in predicting the source of secondary structure in unbound molecules as opposed to bound molecules. Consider the top two ROC curves in Figure 2. In each of the two experiments a classifier was constructed by thresholding the T-S ambiguity index, declaring the secondary structure to be “comparative” when $d_{T-S}(p, s) < t$ and “MFE” otherwise, since, as we have seen, smaller values of the ambiguity indexes are generally associated with the comparative secondary structure. The ROC curves are then swept out by varying t . The difference between the two panels is in the population used for the classification experiments—unbound molecules in the upper-left and bound molecules in the upper-right. In unbound molecules, small values of the T-S ambiguity index are a good indication that a secondary structure was derived from comparative rather than free-energy analysis, but not in bound molecules. The corresponding hypothesis tests seek evidence against the null hypotheses that in a given group (unbound or bound) the set of positive T-S ambiguity indexes ($d_{T-S}(p, s) > 0$) are equally distributed between the comparative and free-energy derived indexes, and in favor of the alternatives that the T-S ambiguity indexes are less typically positive for the comparative secondary structures. The necessary data can be found in Table 3. The results are consistent with the classification experiments: the hypergeometric p-value is $5.4 \cdot 10^{-14}$ for the unbound population and 0.07 for the bound population.

The same experiments, with the same conclusions, were also performed using the D-S ambiguity index, as shown in the bottom two panels of Figure 2, for which the corresponding hypergeometric p-values are $3.8 \cdot 10^{-7}$ (unbound population) and 0.01 (bound population).

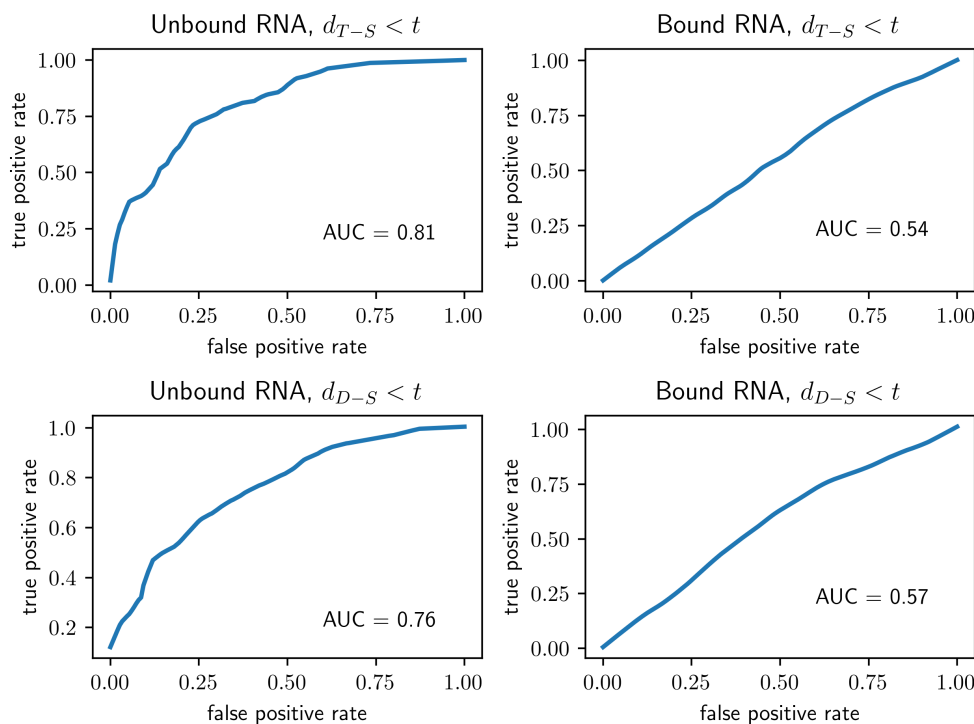


Fig 2. Comparative or MFE? As in Figure 1, each panel depicts the ROC performance of a classifier based on thresholding the T-S (top two panels) or D-S (bottom two panels) ambiguity indexes. Here, small values are taken as evidence for comparative as opposed to MFE secondary structure. Either index, T-S or D-S, can be used to construct a good classifier of the origin of a secondary structure for the unbound molecules in our data set (left two panels) but not for the bound molecules (right two panels). Conditional p-values were also calculated, using the hypergeometric distribution and based only on the signs of the indexes. In each case and the null hypothesis is that comparative secondary structures are as likely to lead to positive ambiguity indexes as are MFE structures, whereas the alternative is that positive ambiguity indexes are more typical when derived from MFE structures: *Upper Left*: $p = 5.4 \times 10^{-14}$; *Upper Right*: $p = 0.07$; *Lower Left*: $p = 3.8 \times 10^{-7}$; *Lower Right*: $p = 0.01$.

References

1. S F Altschul and B W Erickson. Significance of nucleotide sequence alignments: a method for random sequence permutation that preserves dinucleotide and codon usage. *Mol. Biol. Evol.*, 2(6):526–538, November 1985.
2. W M Fitch. Random sequences. *J. Mol. Biol.*, 163(2):171–176, January 1983.
3. D Kandel, Y Matias, R Unger, and P Winkler. Shuffling biological sequences. *Discrete Appl. Math.*, 71(1):171–185, December 1996.
4. Srividya Mohan, Chiaolong Hsiao, Halena VanDeusen, Ryan Gallagher, Eric Krohn, Benson Kalahar, Roger M Wartell, and Loren Dean Williams. Mechanism of RNA double Helix-Propagation at atomic resolution. *J. Phys. Chem. B*, 113(9):2614–2623, March 2009.

5. D Pörschke. A direct measurement of the unzipping rate of a nucleic acid double helix. *Biophys. Chem.*, 2(2):97–101, August 1974.
6. D Pörschke. Model calculations on the kinetics of oligonucleotide double helix coil transitions. evidence for a fast chain sliding reaction. *Biophys. Chem.*, 2(2):83–96, August 1974.
7. D Pörschke. Elementary steps of base recognition and helix-coil transitions in nucleic acids. *Mol. Biol. Biochem. Biophys.*, 24:191–218, 1977.