

1 default

2 An RNA fast-folding path heuristic

3 Vaitea Opuu, Nono S. C. Merleau, and Matteo Smerlak

4 Max Planck Institute for Mathematics in the Sciences, Leipzig, Germany

5 April 10, 2021

6 **Abstract**

7 We propose a heuristic to the folding dynamic making use of a mirror encoding and
8 the fast Fourier transform (FFT) called RAFFT. Based on simple folding rules, it can
9 create many parallel folding paths. The performance in the folding task on a well-curated
10 dataset when compared to the state-of-the-art folding tools was fair. However, when
11 all parallel folding paths were analyzed, it revealed near-native predictions (79% PPV
12 and 81% sensitivity) for sequences of length below 200 nucleotides. On average, those
13 predictions were found to be of similar quality to recent deep-learning-based methods.
14 The folding paths were built with the stem rate model which displays coarse-grained
15 folding paths. Stems are sequentially added during the folding if it improves the overall
16 stability. Those two simple rules create smooth coarse-grained folding paths which are
17 intuitive to analyze and get along with the traditional two states view of the protein
18 folding landscape. Hence, those paths could well approximate fast-folding paths. Since
19 the algorithm was designed toward speed, it can readily be applied to large RNAs.

20 **Introduction**

21 Natural RNA molecules as proteins have biologically relevant functions in many cellular contexts
22 such as protein translation (mRNA, tRNA...), but also in protein-like functions where RNA

23 can perform enzyme functions. Generally, those functions can be better understood through
24 the light of their static tri-dimensional structure. However, some important regulatory RNAs
25 like riboswitch have their biological function tightly bound to their dynamic behavior [24].
26 Therefore, a good understanding of RNAs dynamical aspects is important.

27 RNA molecules are bio-polymers composed of nucleotides. These nucleotides are simple
28 molecules composed of a phosphate-based backbone, a ribose, and a variable nucleobase. Four
29 different canonical bases/nucleobases compose the RNA, namely adenine (A), guanine (G),
30 uracil (U), and cytosine (C). As amino-acid sequences, these nucleotide sequences can form
31 complex tri-dimensional structures critical for their biological functions. Three main structure
32 types are generally considered: the primary structure which is the nucleotide sequence itself.
33 The secondary structure is defined by interacting pairs of nucleobases called base pairs. Next,
34 the tertiary interactions involve other weaker non-trivial interactions within the same sequence.
35 Unlike proteins, RNA structures are usually hierarchically formed. The secondary structure
36 is formed first, followed by the tertiary structure [22]. Moreover, the secondary structure
37 provides an accurate enough description of the thermodynamics and kinetics of RNA molecules.
38 Although base pairs can be formed with various configurations [11], we considered here only
39 the canonical base pairs edge-to-edge interactions: G-C, A-U, and G-U. Many subtleties can
40 be used to define the secondary structure, but we used here a well-accepted formal definition
41 called pseudoknot-free. This forces the RNA secondary structure to be drawable onto a plan
42 where base pairs cannot cross. In the rest of this work, structure refers to the RNA secondary
43 structure.

44 The structure space of an RNA molecule is described by the stability of individual possible
45 structures. The stability ΔG_s of a structure s is the free energy changes with the completely
46 unfolded state. To predict secondary structures, thermodynamic-based methods use empirical
47 data to estimate RNA stability. By assuming the additivity [4] of the loop contributions to
48 the overall stability, the nearest-neighbor loop energy model is the most used model [23]. It
49 is a tabulated set of parameters associating free energy values, determined by optical melting
50 experiments, to individual loop types and compositions such as the well known Turner2004 set

of parameters [16]. Its functional form allows for generalizable energy parameters and the use of an efficient dynamic programming algorithm. It can determine the minimum free energy (MFE) structure of a sequence in the structure space. The MFE is considered a gold standard for free-energy-based predictions. Other estimates exist such as the maximum expected accuracy (MEA), however, it was not found to be significantly better than the MFE [15]. Also, the MFE has an intuitive interpretation. Several tools implement this algorithm, namely Zuker algorithm [26], such as RNAfold [8], Mfold [25], or RNAstructure [18]. Although those methods were found to be consistently accurate at predicting RNA secondary structures as shown in recent benchmarks [20, 9], the additivity foundation is expected to be doomed when sequences get larger and structures complexify. The non-additivity of tertiary interactions and pseudoknots pairings can partially explain this discrepancy. Pseudoknots loop are not defined in the main parameters sets like the Turner2004 model. Another limit of such structure estimates is the static picture that it gives to the RNA folding landscape. From a dynamical standpoint, the sequence navigates the structure space by following the landscape drawn by the stability.

Dynamical information on RNA molecules was found to give valuable complementary information. To follow the dynamic of individual RNA molecules, three rate models describing elementary steps in the structure space are currently used. The base stack model uses base stacks formation and breaking as elementary moves . The base pair rate model uses base pairs as elementary steps as implemented in kinfold [6]. kinfold uses a continuous-time Monte Carlo simulation to follow the RNA folding. It gives the finest resolution in the secondary structure folding landscape, but at the cost of computation time. The stem model uses the creation or deletion of stems to construct the folding dynamics. It is the first strategy explored [14], and provides a coarse-grained description of the kinetic. The folding rates are determined by the free energy changes when stems are added or removed. Although none of these models were definitively rejected nor accepted, this one makes a notable assumption. Indeed, transition states (or saddle points) hidden in the formation of a given stem are not considered. An alternative approach, implemented in kinwalker [7], used the observation that folded intermediates are generally locally optimal conformations. Therefore, locally optimal structures are formed

79 using the standard dynamic programming algorithm and aggregated together along with the
80 folding dynamic.

81 From folding experiments, Pan and coworkers found parallel pathways for a ribozyme which
82 involve two types of path to reach the native structure [17]. One population of sequences was
83 found to fold rapidly, and one quickly reached metastable misfolded structures that slowly fold
84 into the native structure. It is a direct consequence of the ruggedness nature of the RNA folding
85 landscape [21]. Russell and coworkers revealed experimentally the presence of deep channels
86 separated by large energy barriers on the folding landscape which lead to the fast and slow
87 folding paths observed [19].

88 Here we propose a complementary approach by approximating fast-folding paths based on
89 simple folding rules. The basic idea is to use the stem rate model to create multiple parallel
90 folding paths. Here, stems are not allowed to be removed and can be formed only if it improves
91 the stability. It uses a mirror encoding and relies on the fast Fourier transform to speed
92 up the search of stems. This method is inspired by MAFFT [10], a well-known multiple-
93 sequence-alignment tool. The mirror encoding is a simple numerical orthogonal representation
94 of nucleotide sequences. Other similar encodings combined with the FFT were developed for
95 the analysis of DNA [5]. To assess the reliability of the paths predicted, we compared its
96 performance on the folding task for a well-curated dataset, archive II [15]. The algorithm is
97 compared to two estimates: the MFE computed by RNAfold and an ML estimate computed
98 with MxFold2, a recent deep-learning based method [20]. Next, we applied the algorithm to a
99 simple test case, the Coronavirus frameshifting stimulation element [2], where it found closer
100 structures than the MFE.

101 **FFT based folding dynamic heuristic**

102 We now describe the heuristic starting from one sequence S and its associated unfolded structure
103 of length L . We first create a numerical representation of S where each type of nucleotide is

¹⁰⁴ replaced by a unit vector of 4 components:

$$A \rightarrow \begin{pmatrix} 1000 \end{pmatrix} U \rightarrow \begin{pmatrix} 0001 \end{pmatrix} C \rightarrow \begin{pmatrix} 0100 \end{pmatrix} G \rightarrow \begin{pmatrix} 0010 \end{pmatrix} \quad (1)$$

¹⁰⁵ which gives us a $4 \times L$ matrix we call X where each row is a nucleotide type channel. Here,
¹⁰⁶ the first row would be the A channel which we refer to as X^A . Then, we create a second copy
¹⁰⁷ for which we revert the order of the sequence and use the following complementary encoding:

$$\bar{A} \rightarrow \begin{pmatrix} 000w_{AU} \end{pmatrix} \bar{U} \rightarrow \begin{pmatrix} w_{AU}w_{GU}00 \end{pmatrix} \bar{C} \rightarrow \begin{pmatrix} 00w_{GC}0 \end{pmatrix} \bar{G} \rightarrow \begin{pmatrix} 0w_{GC}0w_{GU} \end{pmatrix} \quad (2)$$

¹⁰⁸ Where \bar{A} (respectively $\bar{U}, \bar{C}, \bar{G}$) is the complementary of A (respectively U, C, G). $w_{AU}, w_{GC},$
¹⁰⁹ w_{GU} are tunable parameters for the next step. We call this new complementary copy \bar{X} , the
¹¹⁰ mirror of X .

¹¹¹ Next, for each of the 4 channels, we compute the correlation between X and \bar{X} and by
¹¹² simply summing up the channel correlations, we obtain the correlation between the two copies:

$$cor(k) = (c_{X^A, \bar{X}^A}(k) + c_{X^U, \bar{X}^U}(k) + c_{X^G, \bar{X}^G}(k) + c_{X^C, \bar{X}^C}(k)) / \min(k, 2 \times L - k) \quad (3)$$

¹¹³ where $c_{X^A, \bar{X}^A}(k)$ is the correlation in the A channel between the two copies. $cor(k)$ gives the
¹¹⁴ average number of base pairs for a positional lag k . One channel correlation between copies is
¹¹⁵ given by:

$$c_{X^A, \bar{X}^A}(k) = \sum_{1 \leq i \leq L, 1 \leq i+k \leq M} X^A(i) \times \bar{X}^A(i+k) \quad (4)$$

¹¹⁶ where $X^A(i)$ and $\bar{X}^A(i+k)$ are the A channel of site i and $i+k$. $X^A(i) \times \bar{X}^A(i+k)$ is non
¹¹⁷ zero if sites i and $i+k$ can form a base pair, and will have the value of the chosen weight as
¹¹⁸ described above. Although this requires $O(N^2)$ operations, it can take advantage of the FFT
¹¹⁹ which reduces drastically its complexity to $O(N \log(N))$.

¹²⁰ The large correlation values between the two copies indicate the positional lag at which
¹²¹ the base pair density is high. Therefore, we use a sliding window strategy to search for the

122 longest consecutive base pairs within the positional lag. Since the copies are symmetrical, we
123 only need to slide over one-half of the positional lag. Once the longest base pairs are identified,
124 we simply compute the free energy change when those base pairs are formed. We perform the
125 same search for the n highest correlation lags, which gives us n potential stems. Then, we add
126 to the current structure the base pairs that give the best change of free energy. Free energies
127 were computed using Turner 2004 energy parameters through Vienna RNA package API [12].

128 We are now left with two segments, the interior, and exterior of the group of consecutive base
129 pairs formed. The two exterior fragments are concatenated together. Then, we simply apply
130 recursively the same procedure on the two segments separately in a "Breadth First" fashion to
131 form new consecutive base pairs, until no base pair formation can improve the energy. Hence,
132 it is straightforward to consider pseudoknots by simply concatenating all the fragments left.
133 When multiple stems can be formed in these independent fragments, we combine those possible
134 independent stems and pick the composition that has the best overall stability.

135 The algorithm described so far tends to be stuck in the first local minima found along the
136 folding trajectory. To alleviate this, we propose a stacking procedure where the best trajectories
137 are stored in a stack and evolved in parallel. Hence, it offers the flexibility of overcoming some
138 energy barriers. Once no stem can be formed, the algorithm stops and output the structure
139 with the best energy found among the structures saved in the stack.

140 Application to the folding task

141 To evaluate the relevance of the folding dynamic heuristic, we compared its performance for the
142 folding task. Also, to assess the effect of sequence lengthens on these predictions, we analyzed
143 their performance length-wise. To localize its performance, we compared with two estimates:
144 the MFE structure computed by RNAfold and the ML-based structure computed by MxFold2.
145 RAFFT predictions were performed using non-optimized weights. 50 structures are formed in
146 parallel for each sequence and 100 positional lags were explored at each step for each of the 50
147 structures.

148 Figure 1 shows the performance in predicted positive values (PPV) and sensitivity for the
149 three methods. It shows that the ML method is consistently better than thermodynamic-based
150 methods. Length-wise T-test between the MFE and ML predictions showed that this difference
151 is significant ($p\text{-value} \approx 10^{-12}$) with a substantial improvement of about 10%. Although RAFFT
152 predictions were found to be comparable to MFE predictions, they are significantly less accurate
153 ($p\text{-value} \approx 0.0002$), with a drastic loss of performance for sequences of length greater than 300
154 nucleotides.

155 Among the 50 structures produced by RAFFT, we found on average at least one prediction
156 with 59% of PPV and 63% of sensitivity as shown figure 1. The overall gain of performances
157 is not significantly different from the MFE predictions. However, for the sequences of length
158 below 200 nucleotides, this gain was found to be substantial and significant ($\approx 16\%$ better
159 than the MFE) with $\text{PVV} \approx 79\%$ and sensitivity $\approx 81\%$. The accuracy for those sequences is
160 equivalent to ML performances. For sequence lengths greater than 300 nucleotides, we observed
161 the same drastic loss of accuracy, although we took only the best prediction among the 50 saved
162 configurations for each sequence. We investigated the dependency to the base pair spanning,
163 however, we did not find any significant effect (see supp. mat.).

164 Two regions of lack of performance were observed for all methods. A group of 28 sequences
165 of length shorter than 80 nucleotides have their known structures at on average 9.8 kcal/mol
166 greater than the MFE structures. Some of them involve large unpaired loops such as displayed
167 in figure 2. The second region is around 200 nucleotides in length. The known structure of
168 those sequences also displayed large unpaired regions 2.

169 To investigate the region of the structure space where the thermodynamic model tends to
170 fail, we computed the composition of the known structures. Loop type lengths were computed
171 in percents. Figure 3 shows principal component analysis (PCA) of those compositions. From
172 the PCA, we observed that the known structures are distributed in the structure space toward
173 interior loops. Also, some natural structures, as shown in figure 2, have large unpaired loops.
174 The center of mass in the principal component space is located in between the high-density
175 stacking and interior loops. This shows that the dataset contains many elongated structures.

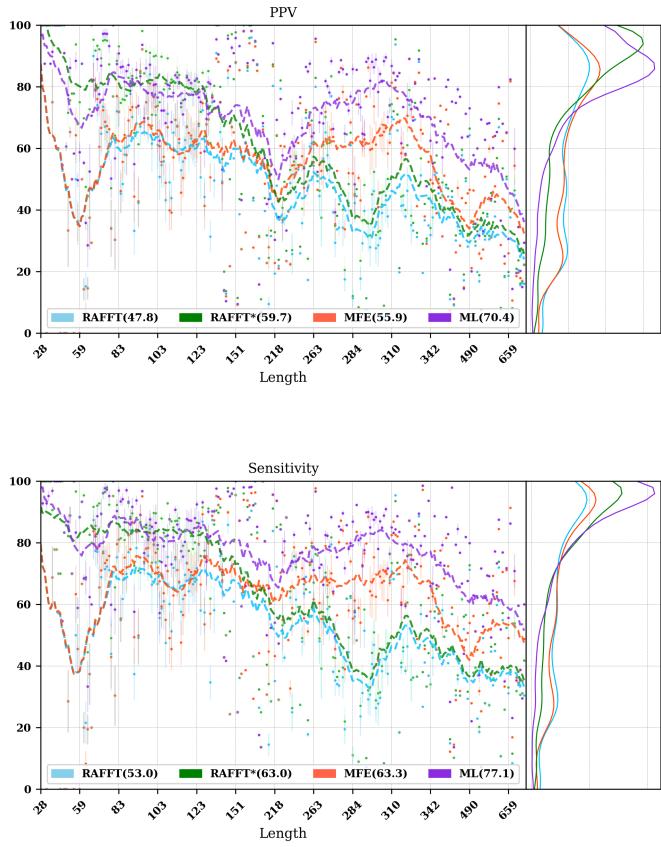


Figure 1: **Predicted positive values and sensitivity results.** RAFFT (blue) displayed the best energy found. RAFFT* shows the best score found among 50 saved structures. Right pанс show the density (sequence-wise) of the accuracy measures.

176 Next, we investigated the structure space produced by the three methods. The thermody-
177 namic approach seems to produce a more diverse structure space as shown in figure 3. Loop
178 contents were extracted from the predicted structures of each method and projected onto their
179 respective two first principal components space. Both RAFFT and MFE predictions seem to
180 produce similar structure spaces while the ML method does allow for long unpaired regions in
181 long hairpins which tend to be closer to the dataset structure space.

182 Test case to predict fast-folding paths

183 Finally, to illustrate RAFFT folding heuristic, we applied it to the Coronavirus frameshifting
184 stimulation element. It is an RNA sequence of about 82 nucleotides with a secondary structure

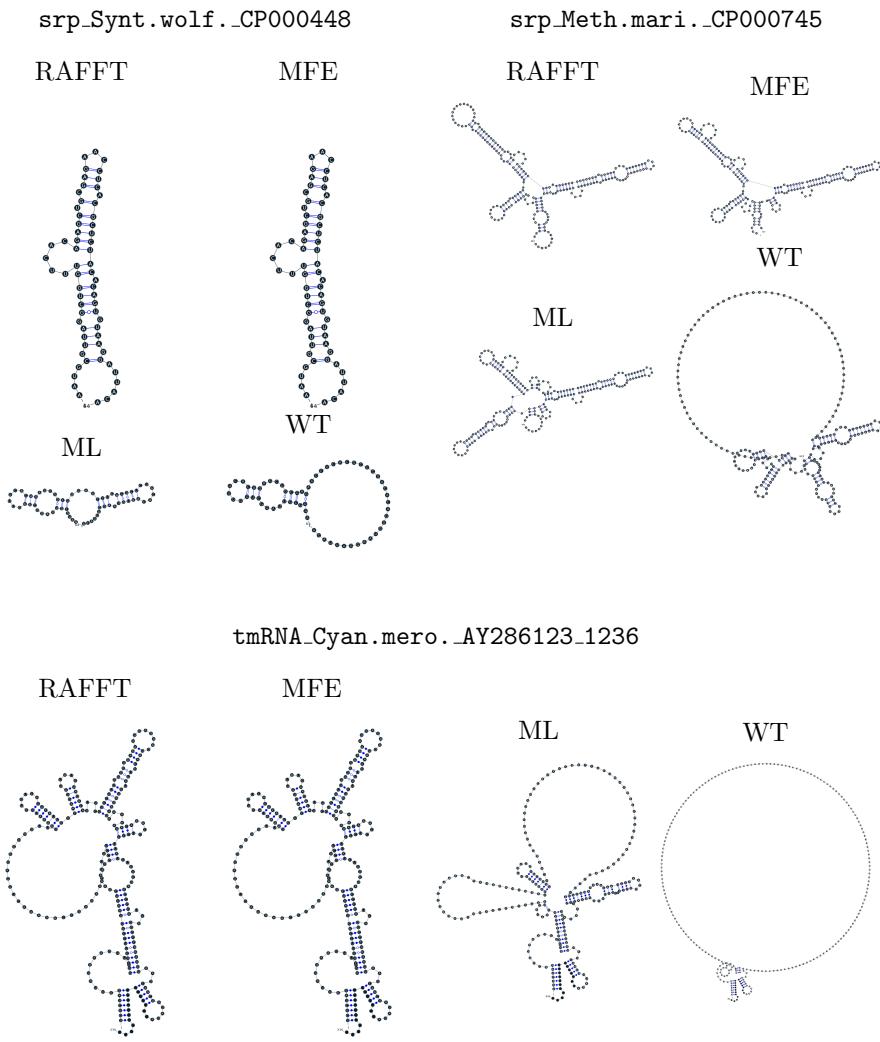


Figure 2: Structures found to be difficult to predict with the thermodynamic model.
The sequence name where extracted directly from the dataset. WT is the known structure.

185 determined by sequence analysis and obtained from the RFAM database. The assumed known
 186 structure has a pseudoknot but was not displayed here. Figure 4 shows the folding path
 187 predicted, the MFE prediction, and the assumed known structure. The approximated fast-
 188 folding paths are predicted in three steps where 5 structures were stored and 100 positional
 189 lags were searched for stems. As shown, some structures explored were not saved or evolved
 190 since no further improvement (relative to all possibilities) was found. RAFFT was able to
 191 recover near-native structures, found to be closer than the MFE, and depicted simple folding
 192 paths. We also tested with 20 saved structure (see supp. mat.), and obtained similar results.

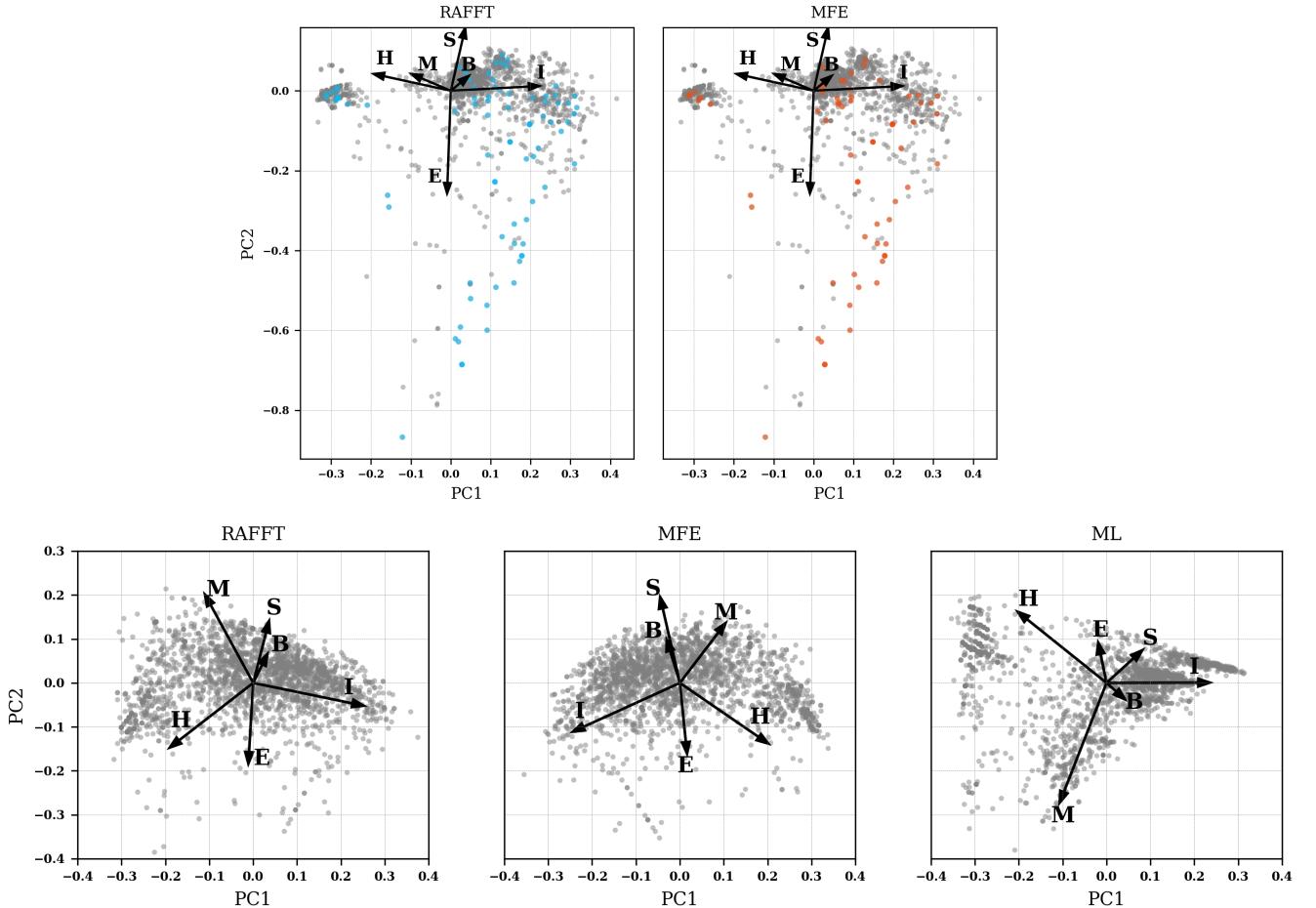


Figure 3: PCA analysis of the known structure and predicted structures. The first row shows two PCAs for the known structures. In the left side, RAFFT predictions with $PPV \leq 10$ are colored in blue. In the left side, MFE structures with $PPV \leq 10$ are colored in orange. The second row shows the PCA of the predicted strucrures for RAFFT, the MFE, and the ML method.

193 However, we observed the greediness effect of the algorithm where a better path in term of
 194 stability. Indeed, a better path was achieve since a more stable structure was found by allowing
 195 less stable intermediates.

196 To visualize the landscape drawn by RAFFT, we produced 300 trajectories with 100 po-
 197 sitional lags explored for stems. All unique structures obtained along each trajectory were
 198 mapped onto a plan using the multidimensional scaling (MDS) algorithm. In the landscape,
 199 the MDS optimized the mapping in such a way that the structure base pair distances were
 200 mostly preserved. Figure 6 shows the landscape interpolated with the 751 unique structures
 201 found. It illustrates the two states folding state where all trajectories started from the high

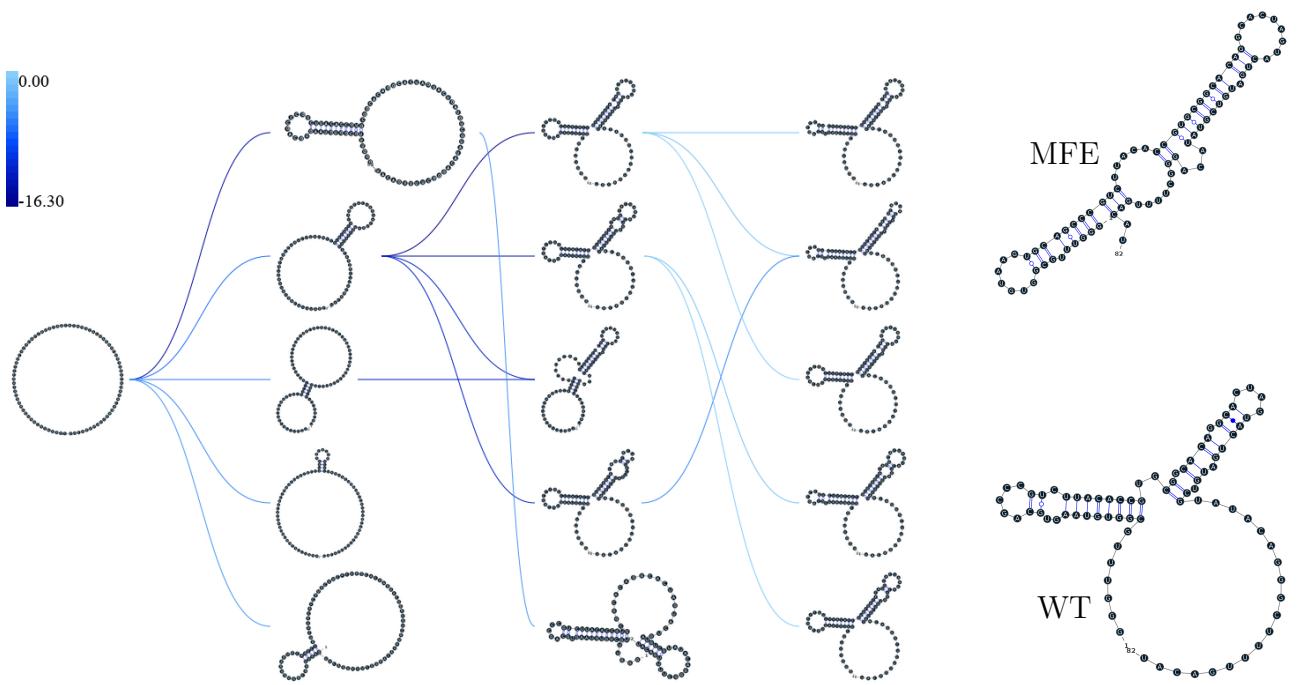


Figure 4: **Fast-folding paths prediction for the Coronavirus frameshifting stimulation element..**

202 peak in the center, and smoothly roll down to the blue area.

203 The graph shown in figure 4 can be used to describe a folding kinetic where transition can
 204 occur from left to right (and right to left) but not vertically. This follows the idea that parallel
 205 paths quickly reach their end points. However, if the end points are non-native states, it will
 206 slowly fold back into the native state. The kinetic is modeled as a Markov process as usually
 207 done [13]. The transition rates $r(x)$ between structures x and y are given by:

$$r(x \rightarrow y) \propto \exp\{-\beta\Delta\Delta G(x \rightarrow y)\} \quad (5)$$

208 where $\beta = 1/kT$ (kcal/mol) is the inverse thermal energy. $\Delta\Delta G(x \rightarrow y)$
 209 is the stability change between structure x and y . Given an initial population of only unfolded
 210 structures, one can simulate the evolution of structure populations. The structure 347 which
 211 dominates the kinetic at the end has a stability of -21.0 kcal/mol, and the MFE (also found by
 212 RAFFT) structure has a stability of 25.8 kcal/mol.

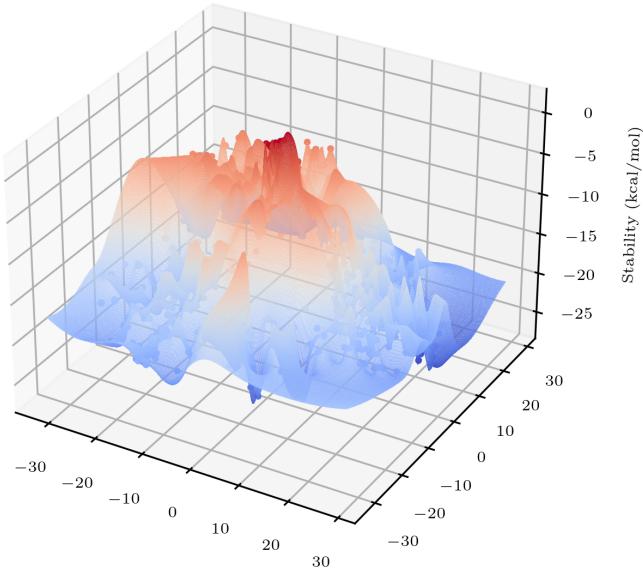


Figure 5: **Landscape built from 300 parallel trajectories.**

213 Concluding discussion

214 We have proposed a heuristic of the RNA fast-fold paths called RAFFT. This heuristic uses
 215 a greedy rules. First, it searches for groups of consecutive base pairs, stems, and from them
 216 if they improve the energy. Hence, it produces smooth and coarse-grained trajectories. To
 217 search for consecutive base pairs, we implemented an FFT-based technique that uses a mirror
 218 encoding. Once a stem is formed, the sequence is split into two independent segments on which
 219 one can recursively search for new stems. For one sequence, the algorithm can follow multiple
 220 folding paths.

221 To assess the relevance of the folding trajectories produced, we compared the algorithm per-
 222 formance for the folding task. Two structure estimates were compared with: the MFE structure
 223 computed using RNAfold, the ML-based estimate using MxFold2. Other thermodynamic-based
 224 and ML-based tools were investigated but not shown here. We chose the MFE since it provides
 225 an intuitive interpretation in the structure landscape, and the MEA prediction was not found

226 to be significantly more accurate [15]. The ML estimates gives a data view of the structure
227 spaces.

228 From our experiments, RAFFT had an overall performance below the MFE predictions by
229 8.1% of PVV and 10.3% of sensitivity. The ML-based approach dominated the predictions
230 (70.4% of PPV and 77.1% of sensitivity). We observed some drastic loss of accuracies when the
231 known structures contained large unpaired regions. However, those sequences were anecdotal in
232 the dataset. Moreover, those regions are unlikely to be stable and assumed to be very flexible.
233 Nevertheless, the effect of unpaired regions seemed less dramatic for the ML method since it
234 can produce some of those atypical structures. No striking evidences of the length effect on
235 prediction quality. In addition, no empirical effects of the base spanning was observed (see
236 supp. mat.) as already pointed out in [1].

237 The PCA performed on the known structure compositions revealed a structure space prone
238 to elongated structures where large unpaired hairpin loops and exterior loops can be observed.
239 The PCA analysis performed on the structures predicted by the thermodynamic-based meth-
240 ods (RAFFT and MFE) shown similar structure space, where unpaired regions are of limited
241 number. On the other hand, the ML method seemed to be closer to the natural structure space.
242 According to the thermodynamic model, those unpaired regions have local stability equal to
243 zero. Hence, those regions are not stable at regular experimental conditions in the sense that
244 they may not have a unique stable structure. However, the ML-method was able to identify
245 such structure more consistently than thermodynamic methods. The PCA revealed a group
246 structures with high percents of hairpins. This may suggest some overfitting effects. Therefore,
247 not being able to recover such structures would be proof of robustness.

248 Although the overall performance of RAFFT was only fair compared in the folding task, we
249 found one among the $k = 50$ predicted trajectories that had better accuracy than the low energy
250 structure displayed. In fact, the gain of performance is substantial for the sequences of length
251 below 200 nucleotides with 16% better in PPV than the MFE predictions. The performance is
252 significantly similar to the ML-base method for that length range. Sequences of length < 200
253 nucleotides represent 86.4% of the total dataset. For the 140 sequences of length greater than

254 300 nucleotides, all k predictions per sequence were similar and performed worst than the other
255 methods. This could be partially explained by the greediness of the algorithm, however, we
256 also believe that the thermodynamic energy model could give a complementary explanation.
257 Indeed, the additivity of the loop contributions to the stability is likely to be doomed for large
258 sequences [22]. However, the MFE did not show any notable discrepancy for large sequences
259 (> 300 nucleotides) except for a few structures with large unpaired regions. This could be
260 explained by the observation used in kinwalker, where locally optimal substructures composed
261 the native structures. Therefore, we assume that the MFE structure is more often composed
262 of locally optimal structures. We tried RAFFT with a larger number of saved structures in the
263 stack, however, it only got closer to the MFE prediction quality and did not perform better
264 (see supp. mat.) on large sequences.

265 As an illustrative example, we applied the heuristic on a natural RNA, the Coronavirus
266 frameshifting element. All trajectories started from the unfolded state to the stable structures
267 in a "two-states" fashion. Furthermore, we showed that the fast-folding paths model can be
268 used as a kinetic model where transitions are given by the parallel paths and intermediate
269 structures found. Because the folding trajectories are already coarse grained and smooth, the
270 kinetic can be drawn without any additional coarsening. Indeed, usual kinetic frameworks
271 ([13]) need a coarse-grained representation of all attraction basins and the saddle points that
272 separate them. The structure dominating the kinetic was a structure close to the native one
273 although lower energy structures where found.

274 Given the experiment results, we believe that RAFFT is a robust heuristic for the fast-
275 folding path since it can produce predictions of high accuracy for 86.4% of this dataset. The
276 folding paths as calculated by RAFFT are smooth and coarse-grained since whole stems are
277 formed, if it improves the energy, and leads to near-native structures. This near-native coarse-
278 grained folding path is an intuitive idea that is similar to the funnel protein folding landscape.
279 We expect this heuristic to give valuable and complementary information to the MFE-like
280 predictions. However, additional efforts are necessary to determine whether the folding paths
281 followed were experimentally observed.

282 On the technical points, the mirror encoding as describe here is a versatile tool for RNA
283 analysis. Since it contains the relative positions of base pairs in the whole sequence, we expect it
284 to be extendable to other use cases such as sequence clustering, or the speed up of Nussinov-like
285 algorithms. On the other hand, we are aware of the limits of choosing the maximal number of
286 base pairs each at each step. However, the greediness of the algorithm had a limited impact on
287 the results. We are not planning to provide yet another folding tool, in this already crowded area
288 of excellent software, but one could combine this tool with an ML-base scoring to discriminate
289 the folding path that is likely to be observed.

290 Methods

291 Starting from the ArchiveII dataset, we first removed all the structures with pseudoknots since
292 all tools considered here don't handle pseudoknots. Next, we removed all the structures which
293 were evaluated with positive or null energy with the Turner 2004 energy parameters. Since
294 positive energies mean that the completely unfolded structure is more stable than the native
295 one. Those structures are assumed not well modeled by the energy function used here and
296 therefore would blur the interpretation of the kinetic we try to extract. This dataset is composed
297 of 2698 structures. 240 sequences were found multiple times (from 2 to 8 times). 19 of them
298 were found with different structures. We discarded all duplication and picked the structure
299 with the lowest energy for each. We obtained a dataset of 2296 sequences.

300 To compute the MFE structure, we used RNAfold (version) with the default parameters
301 and the Turner 2004 set of energy parameters. For the machine learning tool, we computed the
302 prediction using Mxfold2 with the default parameters. Therefore, only one structure prediction
303 per sequence for those two methods were used for the statistics.

304 Two parameters are critical for RAFFT, the number of positional lags in which stems are
305 searched and the number of saved configurations in the stack. For the experiments, we search
306 for stems in the 100 best positional lags and stored 50 conformations. For the predictions
307 analysis, we displayed the lowest energy found at the end for each structure and the most

308 accurate prediction among the 50 saved structures. The correlation which allow to choose the
309 positional lags was computed using the weights $w_{GC}=3$, $w_{AU}=2$, and $w_{GU}=1$.

310 To measure the prediction accuracy, we used two metrics from epidemiology. The positive
311 predictive value (PPV) is the fraction of correct base pairs predictions in the predicted structure.
312 The sensitivity is the fraction of correctly predicted base pairs in the true structure. Both
313 metrics are defined as follow:

$$PPV = \frac{TP}{TP + FN} \quad \text{Sensitivity} = \frac{TP}{TP + FP} \quad (6)$$

314 where TP, FN, and FP stand respectively for the number of correctly predicted base pairs (true
315 positives), the number of base pairs not detected (false negatives), and the number of wrongly
316 predicted base pairs (false positives). To maintain consistency with previous and future studies,
317 we computed these metrics using the implementation in the **scorer** tool provided in [15], which
318 provide also a more flexible estimate where shifts are allowed.

319 The loop compositions were extracted in terms of percent of the cumulative loop sizes.
320 This method, although not accurate, gives an overall idea of the structure space. We first
321 convert the structures into Shapiro notation using Vienna Package API. From the notation, we
322 extracted the sizes of interior, exterior, bulge, stacking, hairpins, and multibranch loops. Next,
323 we converted those sizes into percents of types of loops from which we computed the principal
324 components. For visual conveniences, the structure compositions were projected onto the first
325 two principal components. The composition arrows represent the eigenvectors obtained from
326 the diagonalization of the covariance matrix.

327 The secondary structure representations were obtained with VARNA [3].

328 References

329 [1] AMMAN, F., BERNHART, S. H., DOOSE, G., HOFACKER, I. L., QIN, J., STADLER,
330 P. F., AND WILL, S. *The Trouble with Long-Range Base Pairs in RNA Folding*. Advances

- 331 in Bioinformatics and Computational Biology. Springer International Publishing, 2013,
332 pp. 1–11.
- 333 [2] BARANOV, P. V., HENDERSON, C. M., ANDERSON, C. B., GESTELAND, R. F.,
334 ATKINS, J. F., AND HOWARD, M. T. Programmed ribosomal frameshifting in decoding
335 the sars-cov genome. *Virology* 332, 2 (2005), 498–510.
- 336 [3] DARTY, K., DENISE, A., AND PONTY, Y. Varna: Interactive drawing and editing of the
337 rna secondary structure. *Bioinformatics* 25, 15 (2009), 1974–1975.
- 338 [4] DILL, K. A. Additivity principles in biochemistry. *Journal of Biological Chemistry* 272,
339 2 (1997), 701–704.
- 340 [5] FELSENSTEIN, J., SAWYER, S., AND KOCHIN, R. An efficient method for matching
341 nucleic acid sequences. *Nucleic Acids Research* 10, 1 (1982), 133–139.
- 342 [6] FLAMM, C., FONTANA, W., HOFACKER, I. L., AND SCHUSTER, P. Rna
343 folding at elementary step resolution. *RNA* 6, 3 (2000), 325–338.
- 344 [7] GEIS, M., FLAMM, C., WOLFINGER, M. T., TANZER, A., HOFACKER, I. L., MID-
345 DENDORF, M., MANDL, C., STADLER, P. F., AND THURNER, C. Folding kinetics of
346 large rnas. *Journal of molecular biology* 379, 1 (2008), 160–173.
- 347 [8] HOFACKER, I. L. Vienna rna secondary structure server. *Nucleic Acids Research* 31, 13
348 (2003), 3429–3431.
- 349 [9] HUANG, L., ZHANG, H., DENG, D., ZHAO, K., LIU, K., HENDRIX, D. A., AND
350 MATHEWS, D. H. Linearfold: Linear-time approximate rna folding by 5'-to-3' dynamic
351 programming and beam search. *Bioinformatics* 35, 14 (2019), i295–i304.
- 352 [10] KATOH, K. Mafft: a novel method for rapid multiple sequence alignment based on fast
353 fourier transform. *Nucleic Acids Research* 30, 14 (2002), 3059–3066.
- 354 [11] LEONTIS, N. B., AND WESTHOF, E. Geometric nomenclature and classification of
355 rna base pairs. *RNA* 7, 4 (2001), 499–512.

- 356 [12] LORENZ, R., BERNHART, S. H., ZU SIEDERDISSEN, C. H., TAFER, H., FLAMM, C.,
357 STADLER, P. F., AND HOFACKER, I. L. Viennarna package 2.0. *Algorithms for Molecular*
358 *Biology* 6, 1 (2011), 26.
- 359 [13] LORENZ, R., FLAMM, C., HOFACKER, I., AND STADLER, P. Efficient computation of
360 base-pairing probabilities in multi-strand rna folding. In *Proceedings of the 13th Inter-*
361 *national Joint Conference on Biomedical Engineering Systems and Technologies* (- 2020),
362 p. nil.
- 363 [14] MARTINEZ, H. M. An rna folding rule. *Nucleic Acids Research* 12, 1Part1 (1984), 323–334.
- 364 [15] MATHEWS, D. H. How to benchmark rna secondary structure prediction accuracy. *Meth-*
365 *ods* 162-163, nil (2019), 60–67.
- 366 [16] MATHEWS, D. H., DISNEY, M. D., CHILDS, J. L., SCHROEDER, S. J., ZUKER, M.,
367 AND TURNER, D. H. Incorporating chemical modification constraints into a dynamic pro-
368 gramming algorithm for prediction of rna secondary structure. *Proceedings of the National*
369 *Academy of Sciences* 101, 19 (2004), 7287–7292.
- 370 [17] PAN, J., THIRUMALAI, D., AND WOODSON, S. A. Folding of rna involves parallel
371 pathways. *Journal of Molecular Biology* 273, 1 (1997), 7–13.
- 372 [18] REUTER, J. S., AND MATHEWS, D. H. Rnastructure: Software for rna secondary struc-
373 ture prediction and analysis. *BMC Bioinformatics* 11, 1 (2010), 129.
- 374 [19] RUSSELL, R., ZHUANG, X., BABCOCK, H. P., MILLETT, I. S., DONIACH, S., CHU, S.,
375 AND HERSCHLAG, D. Exploring the folding landscape of a structured rna. *Proceedings of*
376 *the National Academy of Sciences* 99, 1 (2001), 155–160.
- 377 [20] SATO, K., AKIYAMA, M., AND SAKAKIBARA, Y. Rna secondary structure prediction
378 using deep learning with thermodynamic integration, 2020.

- 379 [21] SOLOMATIN, S. V., GREENFELD, M., CHU, S., AND HERSCHLAG, D. Multiple native
380 states reveal persistent ruggedness of an rna folding landscape. *Nature* 463, 7281 (2010),
381 681–684.
- 382 [22] TINOCO, I., AND BUSTAMANTE, C. How rna folds. *Journal of Molecular Biology* 293, 2
383 (1999), 271–281.
- [23] TURNER, D. H., AND MATHEWS, D. H. Nndb: the nearest neighbor parameter database
for predicting stability of nucleic acid secondary structure. *Nucleic Acids Research* 38,
suppl_1(2009), D280 – – D282.
- [24] VITRESCHAK, A. Riboswitches: the oldest mechanism for the regulation of gene expression?
385 *Trends in Genetics* 20, 1 (2004), 44–50.
- [25] ZUKER, M. Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic
387 Acids Research* 31, 13 (2003), 3406–3415.
- [26] ZUKER, M., AND STIEGLER, P. Optimal computer folding of large rna sequences using
389 thermodynamics and auxiliary information. *Nucleic Acids Research* 9, 1 (1981), 133–148.

