

1 From RNA folding to inverse folding:
2 *a computational study*

3 Der Fakultät für Mathematik und Informatik
4 der Universität Leipzig
5 eingereichte

6 D I S S E R T A T I O N

7 zur Erlangung des akademischen Grades
8 DOCTOR RERUM NATURALIUM
9 (Dr.rer.nat.)

10 im Fachgebiet
11 Informatik
12 vorgelegt

13 von Diplominformatiker **Nono Saha Cyrille Merleau**
14 geboren am 26-03-1992 in Bafoussam, Kamerun

15 Leipzig, den December 2021

[August 4, 2022 at 19:02 – 1.0]

¹⁶

Ohana means family.

¹⁷

Family means nobody gets left behind, or forgotten.

¹⁸

— Lilo & Stitch

¹⁹

Dedicated to my loving dad Michel Saha.

[August 4, 2022 at 19:02 – 1.0]

20 ABSTRACT

21 Since the discovery of the structure of Deoxyribonucleic acid
22 ([DNA](#)) in the early 1953s, and its double-chained complement of
23 information hinting at its means of replication, biologists have
24 recognized the strong connection between molecular structure
25 and function. In the past two decades, there has been a surge
26 of research on an ever-growing class of Ribonucleic acid ([RNA](#))
27 molecules that are non-coding but whose various folded struc-
28 tures allow a diverse array of vital functions. From the well-
29 known splicing and modification of ribosomal [RNA](#), non-coding
30 RNAs ([ncRNAs](#)) are now known to be intimately involved in pos-
31 sibly every stage of [DNA](#) translation and protein transcription, as
32 well as [RNA](#) signalling and gene regulation processes.

33 Despite the rapid development and declining cost of mod-
34 ern molecular methods, they typically can only describe [ncRNA](#)'s
35 structural conformations *in vitro*, which differ from their *in vivo*
36 counterparts. Moreover, it is estimated that only a tiny fraction
37 of known [ncRNA](#) has been documented experimentally, often at a
38 high cost. There is thus a growing realization that computational
39 methods must play a central role in the analysis of [ncRNAs](#). Not
40 only do computational approaches hold the promise of rapidly
41 characterizing many [ncRNAs](#) yet to be described, but there is
42 also the hope that by understanding the rules that determine
43 their structure, we will gain better insight into their function
44 and design. Many studies revealed that the [ncRNA](#) functions are
45 performed by high-level structures that often depend on their
46 low-level structures, such as the secondary structure. This thesis
47 studies the computational folding mechanism and inverse folding
48 of [ncRNAs](#) at the secondary level.

49 In this thesis, we describe the development of two bioinfor-
50 matic tools that have the potential to improve our understanding
51 of [RNA](#) secondary structure. These tools are as follows: (1) RAFFT
52 for efficient prediction of pseudoknot-free [RNA](#) folding pathways
53 using the fast Fourier transform ([FFT](#)); (2) aRNAque, an evolution-
54 ary algorithm inspired by Lévy flights for [RNA](#) inverse folding
55 with or without pseudoknot (A secondary structure that often
56 poses difficulties for bio-computational detection).

57 The first tool, RAFFT, implements a novel heuristic to predict
58 [RNA](#) secondary structure formation pathways that has two com-

59 ponents: (i) a folding algorithm and (ii) a kinetic ansatz. When
60 considering the best prediction in the ensemble of 50 secondary
61 structures predicted by RAFFT, its performance matches the recent
62 deep-learning-based structure prediction methods. RAFFT also
63 acts as a folding kinetic ansatz, which we tested on two RNAs: the
64 coronavirus frameshifting stimulation element (CFSE) and a clas-
65 sic bi-stable sequence. In both test cases, fewer structures were
66 required to reproduce the full kinetics, whereas known methods
67 (such as Treekin) required a sample of 20,000 structures and
68 more.

69 The second tool, aRNAque, implements an evolutionary algo-
70 rithm (EA) inspired by the Lévy flight, allowing both local global
71 search, and which supports pseudoknotted target structures. The
72 number of point mutations at every step of aRNAque EA is drawn
73 from a Zipf distribution. Therefore, our proposed method in-
74 creases the diversity of designed RNA sequences and reduces the
75 average number of evaluations of the evolutionary algorithm.
76 The overall performance showed improved empirical results com-
77 pared to existing tools through intensive benchmarks on both
78 pseudoknotted and pseudoknot-free datasets.

79 In conclusion, we highlight some promising extensions of the
80 versatile RAFFT’s method to RNA-RNA interaction studies. We also
81 provide an outlook of both tools’ implications in studying evolu-
82 tionary dynamics.

83 PUBLICATIONS

84 This thesis presents our contributions to RNA secondary struc-
85 tures' computational methods for the folding and inverse folding
86 problems. They were obtained in collaboration with my advi-
87 sor Matteo Smerlak, Vaitea Opuu and Vincent Messow. Most of
88 the ideas and figures have appeared previously in the following
89 publications:

- 90 • [128] **Nono SC Merleau** and Matteo Smerlak (2021). *A*
91 *simple evolutionary algorithm guided by local mutations for an*
92 *efficient RNA design.* In: *Proceedings of the Genetic and Evolu-*
93 *tionary Computation Conference.* pp. 1027-1034.
- 94 • [138] Vaitea Opuu, **Nono SC Merleau**, Vincent Messow,
95 and Matteo Smerlak(2021). *RAFFT: Efficient prediction of*
96 *RNA folding pathways using the fast Fourier transform.* In:
97 *bioRxiv* (**Submitted** to PLoS Comput. Biol.)
- 98 • [129] **Nono SC Merleau** and Matteo Smerlak (2022). *An*
99 *evolutionary algorithm for inverse RNA folding inspired by Lévy*
100 *flights.* In: *bioRxiv* (**Submitted and Accepted**) (At BMC Bioin-
101 *formatics*).

102 In addition to these works in RNA folding and inverse folding, I
103 studied the fragility of RNA viruses during my PhD using multi-
104 agent evolutionary algorithm simulations. I also contributed to
105 various works in natural language processing and multi-agent
106 simulations for Holonification models. None of these investiga-
107 tions,

- 108 • Igor Haman Tchappi, Stéphane Galland, Vivient Corneille
109 Kamla, Jean-Claude Kamgang, **Cyrille Merleau S Nono**,
110 and Hui Zhao (2019). *Holonification model for a multilevel*
111 *agent-based system.* In: *Personal and Ubiquitous Computing*
112 *23(5).*
- 113 • Ivan P Yamshchikov, **Cyrille Merleau Nono Saha**, Igor
114 Samenko, Jürgen Jost (2020). *It Means More if It Sounds*
115 *Good: Yet Another Hypothesis Concerning the Evolution of Poly-*
116 *semous Words.* In: *Proceedings of the 5th International Con-*
117 *ference on Complexity, Future Information Systems and Risk*
118 *(COMPLEXIS 2020), pages 143-148.*

- ¹¹⁹ • **Nono SC Merleau**, Sophie Pénisson, Philip J Gerrish, San-
¹²⁰ tiago F Elena, and Matteo Smerlak (2021). *Why are viral*
¹²¹ *genomes so fragile? The bottleneck hypothesis*. In: *PLoS. Comput*
¹²² *Biol.* 17(7).

¹²³ will be addressed in this manuscript.

We have seen that computer programming is an art, because it applies accumulated knowledge to the world, because it requires skill and ingenuity, and especially because it produces objects of beauty.

— Donald E. Knuth [104]

129 ACKNOWLEDGEMENTS

¹³⁰ Acknowledgements to be put here.

131 Many thanks to everybody who is already reading through
132 this first draft!

[August 4, 2022 at 19:02 – 1.0]

133 CONTENTS

134	1	INTRODUCTION	1
135	1.1	Survey	1
136	1.2	Characteristics and biological functions of ncRNA	3
137	1.3	Recent advancements in determining ncRNA functions	4
138	1.4	Biochemistry of RNA molecules	6
139	1.5	Bioinformatic definitions	10
140	1.5.1	Structural definitions	11
141	1.5.2	Thermodynamic definitions	15
142	1.5.3	Structural distance definitions	18
143	1.5.4	RNA folding map properties	21
144	1.5.5	The fast Fourier transform (FFT) and evolutionary algorithm (EA) applied to RNA bioinformatics	22
145	1.6	Conclusion and outline of the thesis	26

149 I RNA FOLDING

150	2	INTRODUCTION TO RNA FOLDING	31
151	2.1	Stability and prediction of RNA secondary structures	31
152	2.1.1	MFE prediction tools for pseudoknot-free RNA sequences using a score-base method	34
153	2.1.2	machine learning (ML)-based methods	37
154	2.1.3	Prediction tools for pseudoknotted RNA sequences	40
155	2.2	RNA kinetics	42
156	2.3	Conclusion	45

160 3 RAFFT: EFFICIENT PREDICTION OF FAST-FOLDING PATHWAYS OF RNAs

161	47		
162	3.1	Material and Methods	47
163	3.1.1	RAFFT's algorithm description	48
164	3.1.2	Kinetic ansatz	51
165	3.1.3	Benchmark datasets.	52
166	3.1.4	Structure prediction protocols	53
167	3.2	Experimental results	55
168	3.2.1	RAFFT's run time and scalability	55
169	3.2.2	Accuracy of the predicted structural ensemble	58
170	3.2.3	Applications to the RNA kinetics	60

172	3.3 Conclusion	65
173	II RNA DESIGN	
174	4 INTRODUCTION TO RNA DESIGN	69
175	4.1 RNA inverse folding and biotechnological impli-	
176	cations	69
177	4.2 The positive and negative design	70
178	4.3 Objective functions previously used in the context	
179	of Inverse RNA folding	71
180	4.4 A review on existing inverse RNA folding tools.	73
181	4.4.1 Pseudoknot-free RNA inverse folding tools	74
182	4.4.2 Pseudoknotted RNA inverse folding tools	79
183	4.5 Benchmarking the Inverse folding tools	80
184	4.6 Conclusion	82
185	5 AN EVOLUTIONARY ALGORITHM FOR INVERSE FOLDING	
186	INSPIRED BY LÉVY FLIGHTS.	83
187	5.1 Material and methods	83
188	5.1.1 aRNAque's mutation operator	83
189	5.1.2 aRNAque's objection functions	86
190	5.1.3 aRNAque's EA	88
191	5.1.4 Benchmark parameters and protocols	90
192	5.2 Experimental results	92
193	5.2.1 aRNAque's performance on pseudoknot-free	
194	target structures	93
195	5.2.2 aRNAque's performance on pseudoknotted	
196	target structures	97
197	5.2.3 Quality of the designed RNA sequences	101
198	5.2.4 Complexity and CPU time comparison	103
199	5.3 Conclusion	106
200	III GENERAL CONCLUSION AND DISCUSSIONS	
201	6 ADVANTAGES AND LIMITATIONS OF THE PROPOSED METH-	
202	ODS	109
203	6.1 RAFFT: Limitations and future works	109
204	6.2 aRNAque: Limitations and perspectives	112
205	6.3 RAFFT, aRNAque and evolutionary dynamics per-	
206	spectives	115
207	6.4 Conclusion	118
208	7 GENERAL CONCLUSION	121
209	IV APPENDIX	
210	A rafft APPENDICES	125
211	A.1 Kinetic comparison	125

212	A.2	RAFFT example calls	126
213	A.3	RAFFT performance analysis for 200 structures saved.	128
214	A.4	RAFFT performance analysis with various values of minimum energy contribution required for loop formation	128
216	A.5		128
218	B	arnaque APPENDICES	133
219	B.1	aRNAque's GC-content parameters	133
220	B.2	Benchmark on Eterna100 dataset	133
221	B.3	General EA benchmark parameters	133
222	B.4	Other benchmark on Eterna100-V1	135
223	B.5	Tools patching	136
224	B.6	aRNAque example calls	137
225	B.7	Lévy flight vs Local search: designing the structure with the smallest neutral set in the space of all RNA sequences of length 12	137
227	B.8	Continuous and discontinuous transitions in evo- lution	138
229			
230		BIBLIOGRAPHY	141

231 LIST OF FIGURES

- 232 Figure 1.1 **RNA nucleotides.** Adenine and guanine belongs into the chemical class of purine molecules and the uracil and thymine in the class of pyrimidines 7
- 233
- 234
- 235
- 236 Figure 1.2 **RNA base-pair interactions.** (a) and (b) are commonly know as Watson-Crick base-pairs. (c) is the wobble base-pair. Hydrogen bonds are indicated as grey dashed lines. Substitution of bold hydrogen residues with ribose 5-phosphate yields the corresponding nucleotides found in RNA molecules. 8
- 237
- 238
- 239
- 240
- 241
- 242
- 243 Figure 1.3 **Pseudoknot patterns found in the PseudoBase++.** For each pseudoknot patterns, the different rows represent respectively the circular and the dotbracket shape representations. The B-type and cH-type are more complex forms of H-type. The full complexity order is H-type < B-type < cH-type < K-type. 10
- 244
- 245
- 246
- 247
- 248
- 249
- 250
- 251 Figure 1.4 **Different secondary structure representations of a random generated RNA sequence.** The minimum free energy (**MFE**) structure is predicted using RNAfold from the ViennaRNA Package [112]. The representation were then drawn using VARNA [31] 13
- 252
- 253
- 254
- 255
- 256
- 257
- 258 Figure 1.5 **RNA secondary structure loop decomposition.** Each loop is highlighted in blue. 14
- 259
- 260 Figure 1.6 **Base-pair probability matrix** of a transfert RNA (**tRNA**) sequence computed using RNAfold 2.4.13. The **MFE** structure is depicted on the left and the sequence on top. The frequency of the **MFE** structure in the structural ensemble Σ_ϕ is 0.116. The dot plot on the right shows the pair probabilities within the equilibrium ensemble as (72×72) -matrix and is an excellent way to visualize structural alternatives. 18
- 261
- 262
- 263
- 264
- 265
- 266
- 267
- 268
- 269

270 Figure 1.7

271

272

273

274

275 Figure 3.1

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292 Figure 3.2

293

294

295

296

297

298

299

300

301

302

303 Figure 3.3

304

305

306

307

308

309

310

311

Evolutionary algorithm flow diagram. The algorithm initializes a population of candidate solutions and then loops over the three genetic operations until the termination criteria are satisfied. [25](#)

Algorithm execution for one example sequence which requires two steps. (Step

1) From the correlation $cor(k)$, we select one peak which corresponds to a position lag k . Then, we search for the largest stem and form it. Two fragments, "In" (the interior part of the stem) and "Out" (the exterior part of the stem), are left, but only the "Out" may contain a new stem to add. (Step 2) The procedure is called recursively on the "Out" sequence fragment only. The correlation $cor(k)$ between the "Out" fragment and its mirror is then computed and analyzing the k positional lags allows to form a new stem. Finally, no more stem can be formed on the fragment left (colored in blue), so the procedure stops.

[50](#)

Fast folding graph constructed using RAFFT.

In this example, the sequence is folded in two steps. The algorithm starts with the unfolded structure on the left. The $N = 5$ best stems are stored in stack 1. From stack 1, multiple stems formation are considered, but only the $N = 5$ best are stored in stack 2. Structures are ordered (from top to bottom) by energy in each stack. All secondary structure visualizations were obtained using VARNA [\[31\]](#).

[51](#)

Execution time comparisons. For samples of 30 sequences per length, we averaged the execution times of five folding tools. The empirical time complexity $O(L^\eta)$ where η is obtained by non-linear regression. RAFFT denotes the naive algorithm(with only $N = 1$ structure saved per stack), whereas RAFFT(50) denotes the algorithm where 50 structures can be saved per stack.

[56](#)

312 Figure 3.4

313

314

315

316

317

318

319

320

321

322

323 Figure 3.5

324

325

326

327

328

329

330

331

332 Figure 3.6

333

334

335

336

Impact of the number of positional lags n and the stack size N on the runtime complexity. For a corresponding length, we generated 30 random sequences, and averaged their execution times. Solid lines indicate the estimated time complexity $O(L^\eta)$ where η is obtained with a non-linear regression on these average execution times for (A) Different number of positional lags n and (B) Different stack sizes N . [57](#)

RAFFT's performance on folding task. (A) positive predictive value ([PPV](#)) vs sequence length. In the top panel, RAFFT (in light blue) shows the [PPV](#) score distributions when for the structure (out of $N = 50$ predictions) with the lowest free energy, whereas RAFFT* (in blue) shows the best [PPV](#) score in that ensemble. (B) Sensitivity vs sequence length. [60](#)

Structure space analysis. principal components analysis ([PCA](#)) for the predicted structures using RAFFT, RNAfold, MxFold2 compared to the known structures denoted "True". [61](#)

337 Figure 3.7

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

Application of the folding kinetic ansatz on CFSE. (A) Fast-folding graph in four steps and $N = 20$ structures stored in a stack at each step. The edges are coloured according to $\Delta\Delta G$. At each step, the structures are ordered by their free energy from top to bottom. The minimum free energy structure found is at the top left of the graph. A unique ID annotates visited structures in the kinetics. For example, “59” is the ID of the MFE structure. (B) MFE (computed with RNAfold) and the native CFSE structure. (C) The change in structure frequencies over time. The simulation starts with the whole population in the open-chain or unfolded structure (ID o). The native structure (**Nat.1**) is trapped for a long time before the MFE structure (**MFE.1**) takes over the population. (D) Folding landscape derived from the 68 distinct structures predicted using RAFFT. The axes are the components optimized by the MDS algorithm, so the base pair distances are mostly preserved. Observed structures are also annotated using the unique ID. MFE-like structures (**MFE.1**) are at the bottom of the figure, while native-like (**Nat.1**) are at the top. 62

365 Figure 3.8

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

Folding kinetics of CFSE using Treekin.

A) Barrier tree of the CFSE. From a set of 1.5×10^6 sub-optimal structures, 40 local minima were found, connected through saddle points. The tree shows two alternative structures separated by a high barrier with the global minimum (**MFE** structure) on the right side. (B) Folding kinetics with initial population I_1 . Starting from an initial population of I_1 , as the initial frequency decreases, the others increase, and gradually the **MFE** structure is the only one populated. (C) Folding kinetics with initial population I_2 . When starting with a population of I_2 , the native structure (labelled **Nat.1**) is observable, and gets kinetically trapped for a long time due to the high energy barrier separating it from the **MFE** structure. 63

384 Figure 3.9

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

RAFFT vs Treekin: folding kinetics of a bi-stable RNA sequence. (A) Barrier tree for the bi-stable example sequence. The local minima and the corresponding barriers are computed from the complete enumeration of the structure space. The bi-stability is visible on the barrier tree through the two branches separated by a high barrier. (B) Folding kinetics trajectories. The left plot shows the folding dynamics starting from a population with I_1 , and the right size is the kinetics when the population is initialized in structure I_2 . When starting from I_1 , S_A is quickly populated; starting from I_2 , the bi-stability is more apparent. (A') Fast-folding graph using RAFFT. A maximum of $N = 20$ structures are stored in a stack at each step and overall 46 distinct structures are visited. (B') Folding kinetics trajectory obtained from the fast-folding graph (indices are different from the barrier tree indices). The dynamics starts with a population with only unfolded structure, and slowly, S_B is populated and gets trapped for a long time before the MFE structure S_A becomes populated. 64

411 Figure 5.1

412
 413
 414
 415
 416
 417
 418
 419
 420
 421
 422
 423
 424
 425
 426
 427
 428
 429
 430
 431
 432
 433
 434
 435
 436

Binomial vs. Zipf distributions. (A) Samplings Binomial and Zipf distributions for the best binomial mutation rate μ^* (respectively c^* for the best Zipf exponent parameter). Both distributions have a mean of 8.7 point mutations for a sequence of length 88 nucleotides. (B) Tuning of binomial mutation rate parameter. For each $\mu \in [0, 1]$ with a step size of 0.005 and the pseudoknotted target PKB00342 of length 88, 50 sequences were designed using aRNAque. (B) shows the median generations and the success percentage *vs.* the mutation rate (μ). The best mutation rate is $\mu^* = 0.085$ (with a median number of generation 93.5 and a success rate of 92%). (C) Tuning of Levy exponent. Similar to (B), for each $c \in [0, 7]$ with a step size of 0.1 and for the same pseudoknotted target structure, 100 sequences were designed using aRNAque. It shows the median generations and the percentage of success *vs.* the exponent parameter (c). The Zipf exponent distribution that produced the highest success rate and the minimum number of generations is $c^* = 1.4$. **85**

437 Figure 5.2

438
 439
 440
 441
 442
 443
 444
 445
 446
 447
 448
 449
 450
 451

Parameter tuning for both binomial and Lévy mutation schemes. (A) Lévy mutation parameter tuning. Histogram of best exponent parameter (c^*) for a set of 81 target structures with different pseudoknot patterns and various lengths. The most frequent best exponent value is 1. (B) Binomial parameter tuning. Histogram of best mutation rate (μ^*) for the same set of 81 target structures with different pseudoknots and various lengths. The most frequent best parameter is the low mutation rate ($\approx 1/L$). For some structures, the best mutation rate is the high one for different lengths as well. **92**

452 Figure 5.3

453

454

455

456

457

458

459

460

461

462

463

464

465

Lévy mutation vs. Local mutation: performance analysis with respect to the base-pair density. The higher the base-pair density is, the more useful the Lévy mutation scheme to speed up the optimization EA. (A) Distributions of number of generations for the low and high base-pair density targets of the Eterna100 dataset. (B) Percentages of targets with low and high base-pair density for the Eterna100 and PseudBase++. (C) The length distributions of the low and high base-pair density pseudoknot-free and pseudoknotted targets. 93

466 Figure 5.4

467

468

469

470

471

aRNAque's performance on a TRIPOD secondary structure. (A) The tripod target structure. (B) aRNAque's solution using the Turner1999 energy parameter sets. (C) aRNAque's solution using the Turner2004 energy parameter sets. 96

472 Figure 5.5

**Lévy mutation mode vs local mutation
(one-point mutation).** (A) Hamming distance distributions *vs.* target structure lengths. (B) Number of generations distributions for different length groups. In both (A) and (B), lower values indicate better performance. The target structures are solvable in less than 100 generations for both mutation schemes and most length groups. Still, the difference in the number of generations gets more significant as the target lengths increase, except for the two last length groups for which both mutation schemes mostly failed. The highest difference in terms of median number of generations is 150 for target lengths in the range [124–144] (respectively 123, 49, 46, 16, 7, 0, 0 for the length ranges [84–104], [64–84], [104–124], [44–64], [24–44], [144–164], [164–184]). Averaging over all length groups, the median number of generations difference between the Levy mutation and the one point mutation is 48 generations.

495
496 Figure 5.6

**aRNAque vs antaRNA on PseudoBase++ dataset
using both IPknot and HotKnots.** Lower values imply better performance. (A, B) Base pair distance distributions of the designed sequences to the target structure for different pseudoknot types. (C,D) Mean base pair distance against target lengths.

99
503
100

504 Figure 5.7

aRNAque vs antaRNA on PseudoBase++ dataset using IPknot: GC-content analysis. (A) Base-pair distance distributions. (B) GC-content distance distributions. The difference between the targeted GC-content and the actual GC-content values. In (A,B), lower values imply better performance. (C) Number of successes realised by both inverse folding tools. Two values are considered: the up value represent the number targets successfully solved for each GC-content value out of the 266 targets benchmarked; the down values represent the number sequences folding into the targeted secondary structure. [102](#)

519 Figure 5.8

aRNAque vs antaRNA on PseudoBase++ dataset using IPknot: Diversity analysis. The positional entropy distributions plotted against the targeted GC-content values. Higher values imply better performance. [103](#)

524 Figure 5.9

central processing unit (CPU) time: RNAinverse vs. aRNAque. Each bubble corresponds to a target structure in EteRNA100 dataset and, their colours are proportional to the length of the targets. In the legend, MHD stands for Median Hamming distance, and the different markers represent—('o') 100% success for both tools—('+') 100% success for aRNAque and not for RNAinverse—('−') for the case both tools fail to find at least one sequence that folds into the target. Underlying the CPU time difference is the inside plot that shows the CPU time (in seconds) as a target length function. [104](#)

- 538 Figure 5.10

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553 Figure 6.1

554

555

556

557

558

559

560

561

562

563

564

565 Figure A.1

566

567

568

569

570

571

572 Figure A.2

573

574

575

576

577

CPU time analysis using Hotknots: antaRNA vs. aRNAque. Each bubble corresponds to a target structure in PseudoBase++ dataset and, their colours are proportional to the length of the targets. In the legend, BP stands for Median base pair distance, and the different markers represent—('o') 100% success for both tools—('+') 100% success for aRNAque and not for antaRNA—('−') for the case aRNAque's designed sequences are of median base pair distances greater than the one of antaRNA. Underlying the CPU time difference is the inside plot that shows the CPU time (in seconds) with respect to the target length. [105](#)

Lévy mutation vs one-point mutation. For the Eterna100 target structure [CloudBeta] 5 Adjacent Stack Multi-Branch Loop, ten independent runs were performed in which a minimum of 10 sequences were designed per run. (A) Max fitness and mean fitness (inset) over time. (B) Distinct sequences vs. Distinct structures over time. (C) Mean Shannon entropy of the population sequences over time for both binomial and Lévy mutation. (D) The max fitness plotted against the entropy over time. [114](#)

Structure ensemble characterization. The upper part shows the average probability summed over the ensembles of structures predicted per sequence with different methods. The bottom part shows the average positional entropy of structures using the dot-bracket notation. [126](#)

Positive predictive values and sensitivity results. RAFFT (blue) displayed the best energy found. RAFFT*(200) shows the best score found among 200 saved structures. Left panels show the density (sequence-wise) of the accuracy measures. [129](#)

578 Figure A.3

579

580

581

582

583

584

585

586

587 Figure A.4

588

589

590

591 Figure B.1

592

593

594 Figure B.2

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

Predictive performance of RAFFT with various values of minimum energy contribution required for loop formation. Positive values for this parameter causes RAFFT to accept destabilizing loops, therefore being less greedy than per default. The performance of RAFFT was not observed to be positively affected by allowing sub-optimal loop formation. [130](#)

Base pair spanning: It shows the percent of base pairs predicted found in the known structures per number of nucleotides between them. [131](#)

Distribution of number of generations need to solve the target T_1 , for both Lévy and Local mutation schemes. [138](#)

Simulation of an RNA population evolving toward a tRNA (See the figure on the right side) target secondary structure. The target was reached after 933 generations (i.e. $\approx 10^5$ replications). The black line shows the average structure distance of the structures in the population to the target structure. The evolutionary history linking the initial structure to the target structure comprises 23. Each structure is labelled by an integer taken from 0 to 22. To each of them corresponds one horizontal line (in red). The top-level corresponds to the initial structure and the bottom the target structure. At each level, a series of red intervals correspond to the periods when the structure was present in the population, and the green curve represents the transition between structures. Only the time axis has a meaning for the red and green curves. [140](#)

615 LIST OF TABLES

616	Table 3.1	Average performance displayed in terms of PPV and sensitivity. The metrics were first averaged at fixed sequence length, limiting the over-representation of shorter sequences. The first two rows show the average performance for all the sequences for each method. The bottom two rows correspond to the performances for the sequences of length ≤ 200 nucleotides. 59
625	Table 5.1	Summary of performance of aRNAque vs the 7 other algorithms benchmarked on EteRNA100-V1 by Anderson-Lee et al. [3] (using the recent energy parameter sets, the Turner2004) 94
630	Table 5.2	Summary of performance of aRNAque vs the 10 other algorithms benchmarked on the non-EteRNA100 by Anderson-Lee et al. [3] 95
634	Table B.1	Mutation parameters used in aRNAque to control the GC-content values. 133
636	Table B.2	Evolutionary algorithm parameter for each benchmarks. 134
638	Table B.3	Different parameters for the base pair distributions 135
640	Table B.4	Success percentage on Eterna100 datasets for each set of mutation parameters. 136

642 LISTINGS

643	Listing A.1	Command line to run RAFFT executable after installation 126
645	Listing A.2	Command line to run RAFFT executable after installation 127
647	Listing A.3	RAFFT's output results 127

648	Listing B.1	Command line to run aRNAque python script	137
649	Listing B.2	aRNAque's output results	137

650 LIST OF SYMBOLS AND ACRONYMS

651	DNA	Deoxyribonucleic acid
652	RNA	Ribonucleic acid
653	ncRNA	non-coding RNA
654	CFSE	coronavirus frameshifting stimulation element
655	EA	evolutionary algorithm
656	lncRNA	long non-coding RNA
657	sncRNA	short non-coding RNA
658	tRNA	transfert RNA
659	rRNA	ribosomal RNA
660	cRNA	coding RNA
661	mRNA	messenger RNA
662	CRISPR	clustered regularly interspaced short palindromic repeats
663		
664	SELEX	systematic evolution of ligands by exponential enrichment
665		
666	MFE	minimum free energy
667	DP	dynamic programming
668	NMR	nuclear magnetic resonance
669	ML	machine learning
670	DNN	deep neural network
671	SCFG	stochastic context-free grammar
672	SVM	support vector machine
673	CLLM	conditional log-linear model
674	NN	nearest neighbour
675	MEA	maximum expected accuracy
676	GPU	graphics processing unit
677	WC	Watson-Crick

678	PCA	principal components analysis
679	PC	principal component
680	PPV	positive predictive value
681	API	application programming interface
682	FFT	fast Fourier transform
683	SAVE	synthetic attenuated virus engineering
684	NP	non-deterministic polynomial-time
685	CPU	central processing unit
686	NMCS	nested monte carlo search
687	MCTS	monte carlo tree search
688	ED	ensemble defect
689	NED	normalized energy distance
690	MPGA	massively parallel genetic algorithm
691	piRNA	PIWI-interacting RNA
692	PAR	promoter-associated RNA
693	miRNA	microRNA
694	snoRNA	small nucleolar RNA
695	DFT	distcrete Fourier transform
696	IDFT	Inverse discrete Fourier transform
697	NMR	nuclear magnetic resonance

698

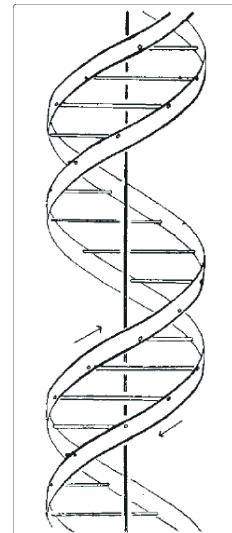
699 INTRODUCTION

700 1.1 SURVEY

701 **DNAs** and **RNAs** are macromolecules in the nucleus of eukaryotic
702 cells that allow storing information with the help of nucleotides.
703 Nucleotides consist of a five-carbon sugar, a phosphate group,
704 and a nucleobase. There are four nucleotides in the **DNA**, distin-
705 guished by their nucleobases: A for Adenine, T for Thymine, G
706 for Guanine, and C for Cytosine. Similar to **DNA**, we also find
707 four different nucleotides in **RNA**, also distinguished by their nu-
708 cleobases with only one exception; the Uracil (U), which replaces
709 Thymine in **DNA**. Even though the basis blocks constituting the
710 **DNA** were known for many years, in 1953, James Watson and
711 Francis Crick [208] succeeded in putting them together and sug-
712 gested a reasonable **DNA** structure. Their work revealed for the
713 first time that the structure of **DNA** molecules has helical chains,
714 each coiled round the same axis where the chain consists of phos-
715 phate dieter groups. The two chains are held together by the
716 purines and pyrimidine bases; they are joined together in pairs,
717 a single base from the other chain bonded to a single base from
718 one chain. For the biding to occur, one of the pairs must be Ade-
719 nine and Thymine or Guanine and Cytosine. A **DNA** molecule
720 structure is depicted on the right side of the page. In contrast to
721 **DNA**, **RNAs** are mostly single-stranded, and the complementary
722 pairings formed in the structure are A-U, G-U and G-C.

723 Watson and Crick's elucidation of **DNA** structure has motivated
724 many other scientists to investigate further the structural implica-
725 tions of molecules in functions such as replication and gave rise
726 to modern molecular biology. Later in the same year, Crick for-
727 mulated the central dogma of molecular biology that describes
728 the flow of information from **DNA** to messenger RNA (**mRNA**)
729 through transcription and from mRNAs to proteins through trans-
730 lation [26]. Since this information flow was proposed, more works
731 have been done to investigate each step.

732 But not all **RNAs** are translated into proteins; in other terms, not
733 all **RNAs** are **mRNAs**. There are mainly two **RNA** groups: coding
734 **RNAs** (**cRNAs**) that are translated into proteins, and non-coding
735 **RNAs** that are not translated into proteins. During the transcrip-



*Helical
representation of
DNA structures
[208].*



The tertiary structure of tRNA.
The CCA-tail is in yellow, the acceptor stem in purple, the variable loop in orange, D-arm in red, the anticodon arm in blue with anticodon in black, and T-arm in green (Taken from Wikipedia)

tion and translation steps in the information flow, some vital functions are performed by ncRNAs such as ribosomal RNA ([rRNA](#)) and [tRNA](#). The tertiary structure of a [tRNA](#) is shown on the left side of the page. The study of such [RNAs](#) revealed that [rRNAs](#), rather than ribosomal proteins, catalyze the synthesis of proteins (i.e. the polymerization of amino acids), distinguish between correct and incorrect codon-anticodon pairs and prevent the premature hydrolysis of peptidyl-tRNAs [15, 133]. Apart from being central to the protein machinery, ncRNAs regulate various biological functions in transcriptional interference, telomere maintenance, epigenetic changes, imprinting, post-transcriptional, translational control, structural organization, cell differentiation and development [51, 159]. We are interested in this work in the structures of ncRNAs.

The function of ncRNAs is largely determined by their high-dimensional structure [19]. For instance, we can analyze the catalytic function of ribozymes in terms of basic structural motifs, e.g. hammerhead or hairpin structures [39]. Other [RNAs](#), like riboswitches, involve changes between alternative structures [202]. Understanding the relation sequence and structure is a central challenge in molecular biology. In the last 20 years, many different methods for determining the [RNA](#) structures of molecules have emerged: from experimental lab methods to computational approaches. For experimental lab methods, X-ray crystallography and the nuclear magnetic resonance ([NMR](#)) are the most accurate approaches to offer structural information at a single base-pair resolution. Both experimental methods are often characterized by high experimental cost and low throughput. In addition to those limitations, [RNA](#) molecules are volatile and difficult to crystallize.

Despite the development of more sophisticated techniques to infer the state of nucleotides in [RNA](#) molecules using enzymatic [95, 200] or chemical probes [193, 212] coupled with next-generation sequencing [12, 192], most of them can only capture [RNA](#) structures *in vitro* which mostly differ from the *in vivo* structure conformations. Experimentally, only a tiny fraction of known ncRNAs has been determined [142]. Because measuring the structure of [RNAs](#) experimentally is very difficult and expensive, computational approaches play a central role in the analysis of natural [RNAs](#) [50, 167], and are an essential alternative to experimental approaches.

Given the ncRNA sequence of bases (primary structure), RNAs fold into secondary structures, such as stem loops and pseudo-knots, before folding into higher level (tertiary and quaternary)

779 structures [16, 194]. This separation of time scales justifies focusing
780 on the secondary structure prediction; evidence suggests that
781 the RNA's secondary structures largely determine the resulting
782 high-level structures [194].

783 This thesis focuses on computational methods addressing RNA
784 molecules' folding and inverse folding at the secondary level.
785 This introductory chapter presents a brief overview of the non-
786 coding RNA concepts. The overview concepts contain biological
787 and biochemical structure definitions of the non-coding RNAs. It
788 also gives an overview of different techniques used to identify
789 new ncRNAs and some applications. It concludes by providing
790 the bioinformatic definitions of RNA secondary structure that
791 constitute the basis and understanding of computational methods
792 and the results presented in this thesis.

793 1.2 CHARACTERISTICS AND BIOLOGICAL FUNCTIONS OF NCRNA

794 In the previous section, we introduced the classical view of information
795 flow in microbiology. Two important ncRNAs involved in
796 the protein machinery have been highlighted (tRNA and rRNA).
797 In this section, we provide some of the main characteristics of
798 ncRNAs, and we emphasize how those characteristics often play
799 an essential role in realizing their functions.

800 What motivates the computational studies of ncRNAs is often
801 the importance of the biological function they play. Consequently,
802 the ncRNAs can be classified based on their biological functions.
803 Although many recent transcriptomic and bioinformatic studies
804 suggested thousands of ncRNAs with their functional importances,
805 the total number of ncRNAs encoded in the human genome still
806 remains unknown [159]. More recently, newly identified ncRNAs
807 have not been validated by their function; it could be possible that
808 most of them are non-functional. Some evolutionary experiments
809 *in vitro* have shown that RNA molecules can catalyze various
810 chemical reactions relevant to biological processes such as RNA
811 replication, nucleotide synthesis, thymidylate synthesis, lipid
812 synthesis, and sugar metabolism [45, 153]. Another characteristic
813 of ncRNAs is their lengths formed post-transcriptionally. We often
814 distinguish two main ncRNA classes of critical biological functions:
815 the short non-coding RNAs (sncRNAs) (with length < 30nt) and
816 the long non-coding RNAs (lncRNAs) (with length > 200nt). The
817 length limit is often because of the practical considerations, in-
818 cluding separating RNAs in standard experimental protocols. The
819 length of ncRNAs is also taken into account in computational stud-

ies, and it will be used throughout our work to distinguish RNA sequences and structures in the different datasets considered.

The function of lncRNAs includes a role in higher-order chromosomal dynamics, telomere biology, and subcellular structural organization [11, 27]. Some lncRNAs play key regulatory and functional roles in the gene expression program of the cell. One of the vital functions is to act as ribozymes. Examples of naturally occurring ribozymes include group I and group II introns—RNase P and the hammerhead. The group, I and group II introns are usually 200–600nt long, catalyzing RNA splicing [65]. Many sncRNAs also contribute to the realization of similar biological functions. For example, small interfering RNAs contribute to gene regulation, transposon control and vital defence. microRNAs (miRNAs) participate in the post-transcriptional gene regulation, miRNAs, PIWI-interacting RNAs (piRNAs)) and promoter-associated RNAs (PARs) contribute to the gene regulation. More recently, many discoveries revealed several ncRNAs implicated in cancer growth and MCL-1 expression regulation [159, 205]. Those examples include ncRNAs from different classes, miRNAs, small nucleolar RNAs (snoRNAs) and T-UCR, all associated with a specific disease [49, 159].

There are also other classes of ncRNAs such as aptamers and riboswitches that have also been observed in nature. Aptamers are ncRNAs that can bind to other specified targets, whose nature is highly diverse. They range from small molecules to larger molecules. In some contexts, aptamers are termed riboswitches; for example, when their function is to sense the presence of an associated metabolite to cause a specific cis-reaction and/or cis-regulation of subordinated functional pathways [213].

In sum, lnc/snc-RNAs contribute to the realization of various biological functions, and they are mostly distinguishable based on their length and functions. But, their functions allow us to distinguish them better. In the next section, we provide some of the recent advancements in the techniques used to identify functional ncRNAs.

1.3 RECENT ADVANCEMENTS IN DETERMINING NCRNA FUNCTIONS

Most of the previously mentioned functions of ncRNAs are identified using gene targeting techniques, a well-known set of techniques used to investigate protein functions [163]. In addition, experimental approaches are used to define ncRNA functions. With the recent advancements in genome engineering, a method

such as clustered regularly interspaced short palindromic repeats (**CRISPR**) has been employed to tag **lncRNAs**, allowing to capture specific **RNA**-protein complexes assembled *in vivo*. This section aims at providing an overview of different techniques used to determine **ncRNA** functions.

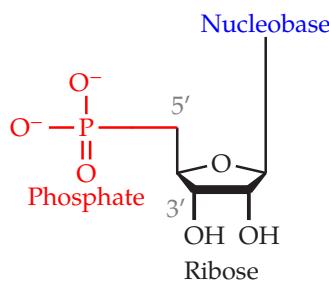
The **CRISPR** [8] was described by Barrangou and his collaborators in 2007 as a distinctive genome feature of most bacteria and archaea and thought to be involved in resistance to bacteriophages. It is an adaptive defence system against viruses and plasmid intrusions. When a successful defence takes place, the system updates information about the intruder's genetic material. This update will then allow the system's host to identify its enemy, making it robust and durable in the future. The information about the intruder's genetic material is stored in short repeating stretches of **RNA**, which can, in the case of a new intrusion, be incorporated into a carrier protein(CAS). The capacities of the **CRISPR/CAS9** of selectively destroying foreign **DNA/RNA** and editing the genome was identified by Li et al. [108], and it was turned into methods allowing to alter and edit single genes within genomes selectively. The same technology is also successfully applied to animal cell lines [85, 91, 204] and industrial plants [109, 186].

Another method of systematic evolution of ligands by exponential enrichment (**SELEX**) [196] introduced by Tuerk in the early 1990s offers the possibility of enriching stretches of **RNA** that can bind a certain target. The method relies on mechanisms usually ascribed to the process of evolution, that is, variation, selection, and replication. A pool of **RNAs** that are entirely randomized at specific positions is subjected to selection for binding, in this case to GP43 on nitrocellulose filters. The selected **RNAs** are amplified as double-stranded **DNA** competent for subsequent *in vitro* transcription. This newly transcribed **RNA** is enriched for better binding sequences and is then subjected to selection to begin the next cycle. Multiple rounds of enrichment result in the exponential increase of the best binding ligands until they dominate the population of sequences. **SELEX** has given rise to numerous synthetic aptamers with different targets in its application. They have been subject to a further extension towards inclusion into regulative **RNA** entities.

More recently, increased types of **ncRNAs** have been detected and identified by the development of next-generation sequencing [206], which can be roughly divided into the process sections

903 of sample preprocessing, library preparation, sequencing, and
904 bioinformatics.

905 The functions of many ncRNAs are dependent on their high-
906 level structures, which often depend on lower-levels sequence
907 and secondary structures. Knowing the structure of an ncRNA
908 plays a vital role in probing its function. For example, Peter Flor
909 and his collaborators used structure information to interpret ex-
910 periments related to the mechanism of RNA function [56]. Or,
911 Yoon et al. suggested new experiments based on RNA secondary
912 structure in yeast to probe RNA functions [96]. Therefore, under-
913 standing even the secondary structure alone can assist both of
914 these examples. In the following section, we provide a biochemical
915 definition of the elementary building blocks of ncRNAs, which
916 are the nucleotides A, U, G and C. In addition, we will provide an
917 overview of the different nucleotide interactions involved during
918 the formation of their secondary structures.



Structure of an RNA nucleotide

919 1.4 BIOCHEMISTRY OF RNA MOLECULES

920 So far, we have provided a biological motivation for studying
921 non-coding RNA as an independent entity. The discovery of new
922 ncRNAs functions has emerged through intensive experimental
923 studies and with recent advanced techniques in next-generation
924 sequencing. Several examples demonstrated the importance of
925 the ncRNA structures in the probing process of new functions.
926 The process in which RNA sequences are mapped to their corre-
927 sponding structures is called RNA folding. In nature, this process
928 is thought to be hierarchical [16, 194]. Nucleotides form a chain
929 given their sequence of bases (primary structure); RNAs fold into
930 secondary structures, such as stem-loops and helices, before fold-
931 ing into higher-level (tertiary and quaternary) structures. Our
932 work is restricted here to the secondary level of an RNA structure,
933 i.e., the set of canonical pairs. This section provides a biochemical
934 definition of different nucleotides and base-pair interactions
935 involved in the secondary structure folding of RNA molecules.

936 Chemically, each nucleotide in RNA molecules consists of a
937 phosphate residue, a pentose sugar and a nucleobase. The typical
938 chemical structure of a nucleotide is depicted on the right side of
939 the page. Figure 1.1 illustrates the chemical structure of each of
940 the four different nucleobases found in RNA (A, C, G and U). A
941 nucleotide is a nucleoside which has a (mono, di, trip) phosphate
942 residue bound to its 5'-carbon atom. By convention, the carbon

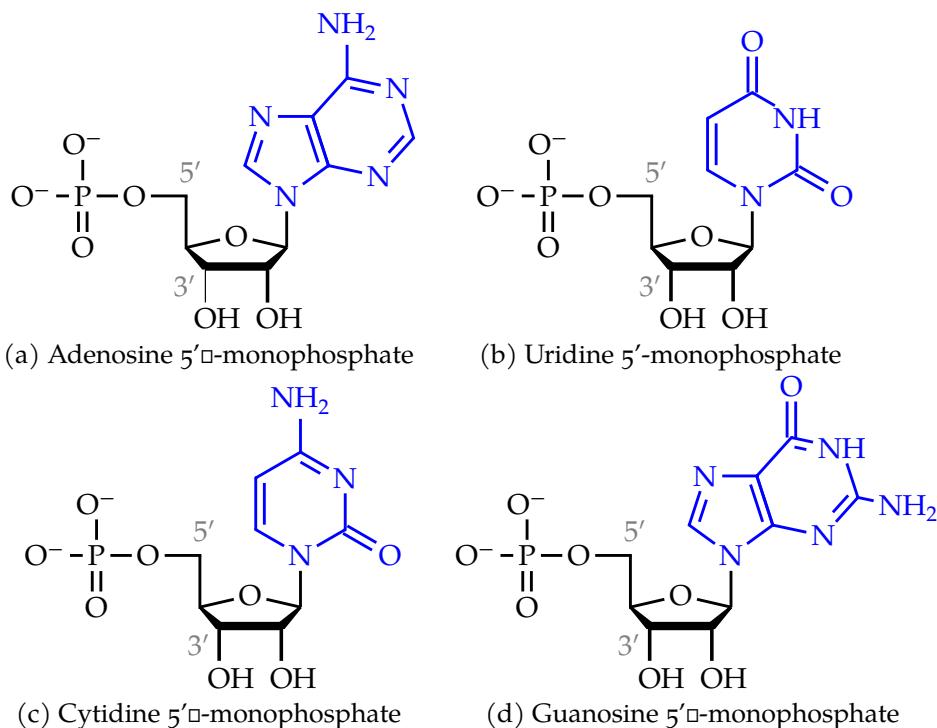


Figure 1.1: **RNA nucleotides.** Adenine and guanine belongs into the chemical class of purine molecules and the uracil and thymine in the class of pyrimidines

943 atoms of the pentose sugar in nucleotides are numbered with
944 *primes*.

945 At the lowest level, **RNA** molecules are simply represented as
946 a list of nucleobase characters. The 5'-3' phosphodiester bonds
947 attach the different nucleotides composing the **RNA** molecule
948 between ribose to form the primary structure of **RNA**. The chain di-
949 rection is conventionally designed as 5' to 3' (i.e. from 5'-phosphate
950 first sugar backbone to the 3'-hydroxyl last sugar in the sequence).

951 In contrast to the **RNA** primary structure, the secondary struc-
952 ture consists of a list of nucleobase pairs, and the hydrogen bonds
953 between the bases form base-pairs. Different interactions are pos-
954 sible between the bases depending on the structure level con-
955 sidered. At the secondary level, we have the Watson-Crick (or
956 canonical) pairs [154, 166] (A-U and G-C), the Wobble (or non-
957 canonical) (G-U) pairs that occur with reduced frequency. **Fig-**
958 **ure 1.2** shows the chemical base-pairs for the Watson-Crick and
959 Wobble interactions.

960 In addition to the Watson-Crick (**WC**) and wobble interactions,
961 we also find crossing or pseudoknotted interactions in natural
962 **RNA** that play vital roles in realizing biological functions, e.g. ribo-

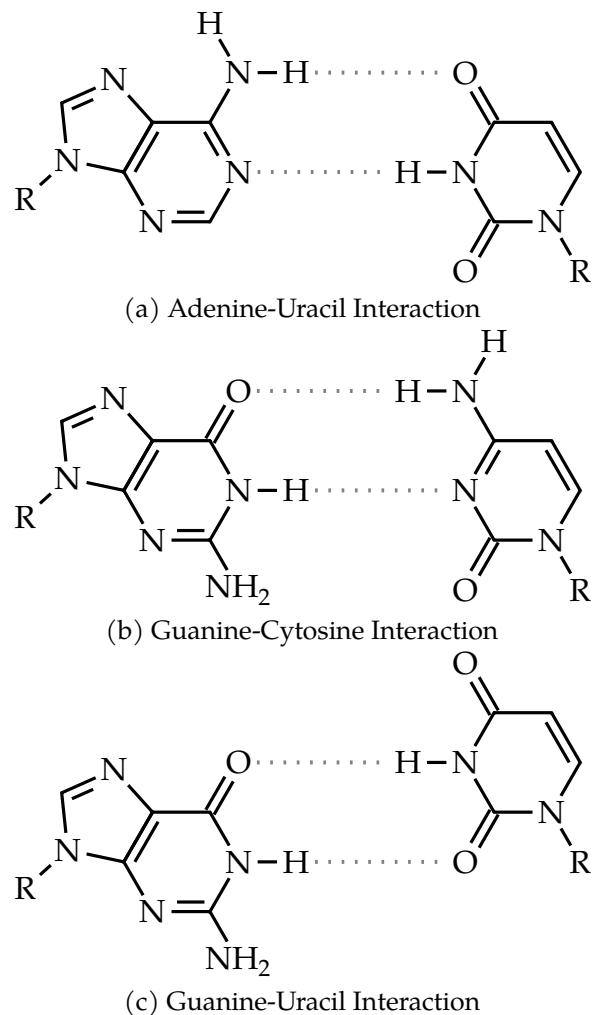
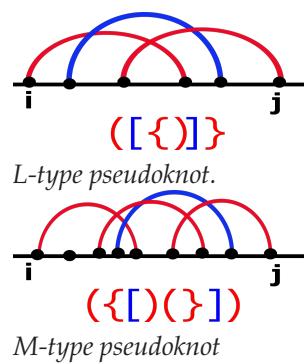


Figure 1.2: **RNA base-pair interactions.** (a) and (b) are commonly known as Watson-Crick base-pairs. (c) is the wobble base-pair. Hydrogen bonds are indicated as grey dashed lines. Substitution of bold hydrogen residues with ribose 5-phosphate yields the corresponding nucleotides found in **RNA** molecules.

somal frame-shifting [63], regulation of translation and splicing [42], or the binding of small molecules [64, 98, 182]. Although the pseudoknots are not considered in the computational folding tool we propose in Chapter 3, they are essential to evaluate the performance of the computational tool we will introduce in Chapter 5 for RNA design. This section also presents different pseudoknot patterns found in natural RNA and emphasizes the one considered in our work.

Pseudoknots occur when two WC, wobble or non-canonical interactions cross each other [210]. Even though pseudoknots are often considered the beginning of the interaction between the secondary and tertiary levels of RNA structures, they account in this work as part of the secondary structure. The restriction to only crossing WC and wobble base pairs contrasts the other tertiary interactions, which may include a broader class of interactions. Many pseudoknot patterns have been identified in natural RNAs. Most occurring pseudoknot patterns tend to be relatively simple in the sense that their crossings are not interlaced and may be viewed as superpositions of two nested secondary structures (bi-secondary structures) [76]. The simplest pattern is often termed as Hairpin-type or H-type (see Figure 1.3a). More complex forms of H-type pseudoknot are bulge hairpin (B-type) or complex hairpin (cH-type). H-type and K-type pseudoknots are the most frequent pseudoknots, but more complex but less frequent pseudoknots are possible such as M-type and L-type pseudoknots (See the figure on the right side of the page) [107]. The four types of pseudoknot patterns considered in Chapter 5 are depicted in Figure 1.3.



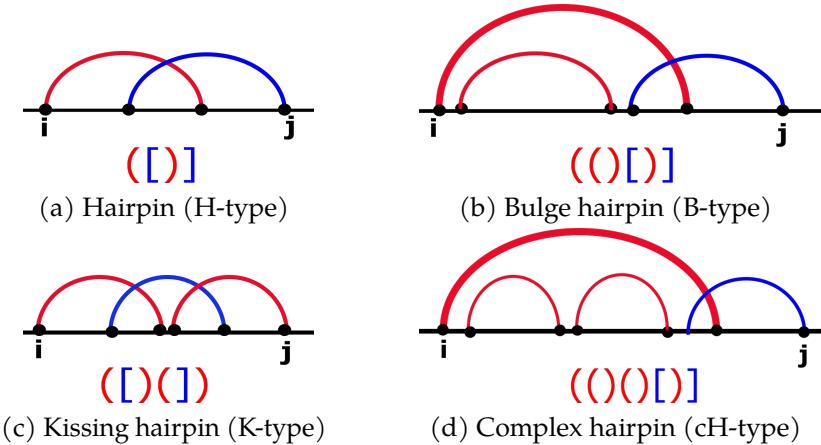


Figure 1.3: Pseudoknot patterns found in the PseudoBase++. For each pseudoknot patterns, the different rows represent respectively the circular and the dotbracket shape representations. The B-type and cH-type are more complex forms of H-type. The full complexity order is H-type < B-type < cH-type < K-type.

Considering pseudoknots in designing functional RNAs is vital given their role in realizing biological functions. Nevertheless, computationally folding an RNA molecule with arbitrary pseudoknot patterns is non-deterministic polynomial-time (**NP**)-complete [118]. Solving this problem is a prerequisite for RNA design and is still a real challenge, not only because of the computational constraint but also the experimental energy measurements of the pseudoknot interactions. In most cases, existing computational tools are restricted to a specific pseudoknot pattern and are based on approximated energy parameters [67].

In the context of this work, we consider two main secondary structure definitions: a pseudoknot-free one in which only canonical interactions with no crossing pairs are allowed and a second one where canonical interactions with possible crossing pairs are permitted. The following section will provide formal definitions and the framework in which we can computationally study the folding of the secondary structure of ncRNAs.

1008 1.5 BIOINFORMATIC DEFINITIONS

We provided in the previous sections the biological motivations and biochemical concepts that support the computation methods studied in the thesis. In order to computationally study and analyse RNA molecules, a more formal representation of RNAs and bioinformatic definitions are required. We provide in this sec-

tion formal definitions and concepts that will support the result presented in this thesis.

1.5.1 Structural definitions

This thesis focuses on computational folding and inverse folding methods of the secondary structure of RNA molecules. The secondary structure, in most cases, is computed for a given RNA sequence. Along the thesis, ϕ will represent an RNA sequence of a fixed length L and \mathcal{S} its corresponding structure. This subsection provides formal definitions of ϕ , \mathcal{S} and the structural properties of \mathcal{S} . We will assume the same definitions in the different tools reviewed in Chapter 2, Chapter 4, which also supports the results presented in Chapter 3 and Chapter 5.

Definition 1 (RNA sequence). More formally, ϕ consists of an ordered sequence of nucleotides that can be represented as:

$$\phi = (\phi_1, \dots, \phi_L), \quad (1.1)$$

where $\phi_i \in \{A, C, G, U\}$ for $i \in \{1 \dots L\}$. ϕ is often known as the primary structure of RNA.

Definition 2 (RNA pseudoknot-free secondary structure). Given an RNA sequence $\phi \in \{A, C, G, U\}^L$, let $\mathcal{P} = \{(i, j) : i < j\}$ be the list of possible pairing positions over the sequence ϕ . A pseudoknot-free secondary structure $\mathcal{S} \subset \mathcal{P}$ of such sequence ϕ is a list of base-pairs with the following constraints [79, 81]:

1. A nucleotide (sequence position) can only belong to a single pair, i.e. $\forall (i, j), (k, l) \in \mathcal{S}$ with $i < k : i = k \Rightarrow j = l$.
2. Paired bases must be separated by at least three unpaired nucleotides. i.e. $\forall (i, j) \in \mathcal{S} \Rightarrow j - i \geq 3$.
3. There are no pseudoknots, i.e. $\nexists (i, j), (k, l) \in \mathcal{S}$ with $i < k < j < l$,
4. The base-pairs consist exclusively of Watson–Crick (C–G and A–U) pairs and Wobble (G–U) pairs. i.e. $\forall (i, j) \in \mathcal{S} \Rightarrow \phi_i \phi_j \in \{GC, CG, AU, UA, GU, UG\}$,

Therefore, RNA secondary structures can be thought as planar graphs than can be more or less easily drawn on a plan.

Definition 3 (Secondary structure representation). A graphical way of representing an RNA secondary structure. There are several representations of \mathcal{S} .

- 1047 • Dot-bracket (or string) representation: In this representa-
 1048 tion, the secondary structure \mathcal{S} is compactly stored in a
 1049 string σ consisting of dots and matching brackets. i.e. σ is
 1050 a string of length L over the alphabet $\Delta_\sigma = \{(.,), [., \{., \}, <$
 1051 , >, .\} where, at each unpaired positions we have a dot '.' at
 1052 the corresponding string position, and $\forall (i, j) \in \mathcal{S}$, we have
 1053 an opening bracket at position σ_i and a closing bracket at
 1054 position σ_j . We denote σ the string representation of the
 1055 structure \mathcal{S} . [Figure 1.4D](#) shows an example of a string rep-
 1056 resentation.
- 1057 • Planar representation: it is the common way of representing
 1058 an RNA secondary structure in which \mathcal{S} is presented as a
 1059 graph with each vertex representing a nucleotide and an
 1060 edge connecting consecutive nucleotides and base-pairs
 1061 (See [Figure 1.4B](#)).
- 1062 • Circular (or circle) representation: similar to planar repre-
 1063 sentation, \mathcal{S} is a graph but drawn in the plane in such a way
 1064 that all vertices are arranged on a circle, and the edges repre-
 1065 senting base-pairs lie inside the circle. In a pseudoknot-free
 1066 secondary structure circular representation, the edges do
 1067 not intersect (See [Figure 1.4A](#)).
- 1068 • Linear representation: In this representation, \mathcal{S} is a graph
 1069 in which the nucleotides are arranged consecutively in a
 1070 line and the edges representing base-pairs form semi-circle
 1071 that do not intersect for pseudoknot-free structure (See
 1072 [Figure 1.4C](#)).
- 1073 • Mountain representation: it is mainly used for representing
 1074 large structures. \mathcal{S} is presented in a two-dimensional graph,
 1075 in which the x -coordinate is the position i of the nucleotide
 1076 in the sequence ϕ and the y -coordinate the number $m(i)$ of
 1077 base-pairs that enclose nucleotide i .
- 1078 • Tree representation: \mathcal{S} is drawn as a tree in which internal
 1079 nodes are the base pairing positions, and the leaves are the
 1080 unpaired positions. The dot-bracket representation is also
 1081 often considered as a tree represented by a string of paren-
 1082 thesis (base-pairs) and dots for the leaf nodes (unpaired
 1083 nucleotides).
- 1084 • Shapiro representation: it allows representing the different
 1085 elements composing \mathcal{S} by single matching brackets, and

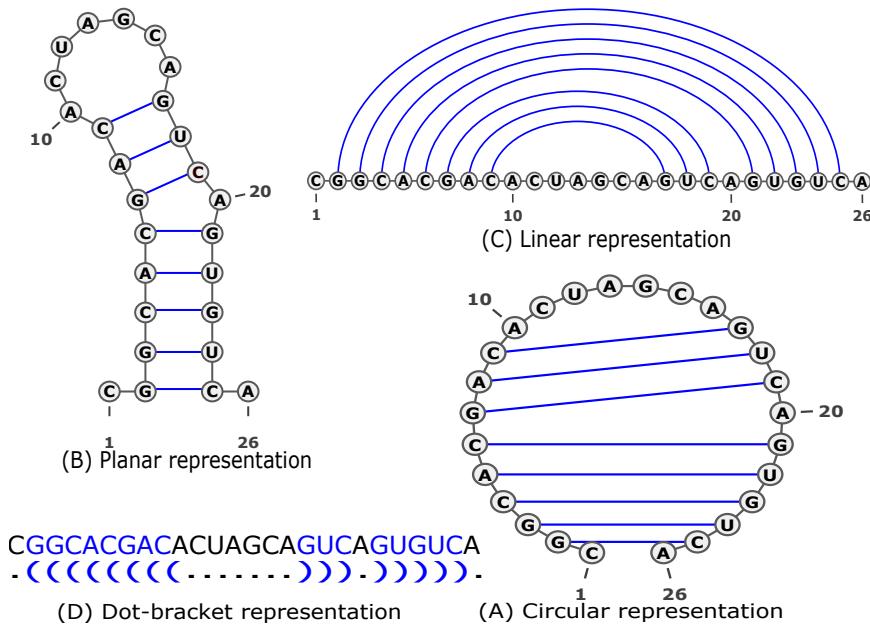


Figure 1.4: Different secondary structure representations of a random generated RNA sequence. The MFE structure is predicted using RNAfold from the ViennaRNA Package [112]. The representations were then drawn using VARNA [31]

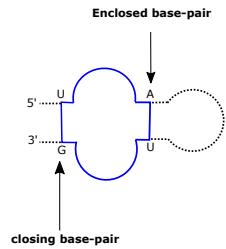
1086 the components are labelled with H(Hairpin), B(Bulge), I
 1087 (interior loop), M (multi-loop) and S (stacking loop) [169].

1088 **Figure 1.4** shows some examples of RNA secondary structure rep-
 1089 resentation. For graphical illustrating examples in the thesis, we
 1090 will mostly use the planar representation, and for computational
 1091 methods, we will use the dot-bracket representation for simplicity.
 1092

1093 **Definition 4** (Secondary structure loop). There exists a unique
 1094 decomposition of \mathcal{S} into a set of n loops $\mathbb{L}_{\phi, \mathcal{S}}$, where loops are the
 1095 faces of its planar drawing. Each loop $\mathcal{L} \in \mathbb{L}_{\phi, \mathcal{S}}$ is characterised
 1096 by its length l (the number of unpaired nucleotides in the loop)
 1097 and its degree d (the number of base-pairs delimiting the loop,
 1098 including the closing loop pair).

1099 By definition, $\forall \mathcal{L} \in \mathbb{L}_{\phi, \mathcal{S}} \Rightarrow \mathcal{L} = \mathcal{L}_p \cup \mathcal{L}_u$ where \mathcal{L}_p and \mathcal{L}_u
 1100 denote respectively the set of loop base-pairs and the unpaired
 1101 positions. \mathcal{L}_p contains only one closing loop and the rest are
 1102 enclosed base-pairs. We say $(i, j) \in \mathcal{L}_p$ is a closing pair if and
 1103 only if $\forall \mathcal{L}_p \ni (i', j') \neq (i, j) : i < i' < j' < j$.

- 1104 1. Interior loop: a loop with degree $d = 2$ i.e $|\mathcal{L}_p| = 2$ and
 1105 $\mathcal{L}_u \subset \{1, 2, \dots, L\} \cup \emptyset$.



- 1106 2. Stacking pair: an interior loop of length $l = 0$ i.e. $|\mathcal{L}_p| = 2$
 1107 and $\mathcal{L}_u = \emptyset$.
- 1108 3. Hairpin Loop: Any loop of degree $d = 1$ and length $l \geq 3$.
 1109 i.e $|\mathcal{L}_p| = 1$ and $\mathcal{L}_u \neq \emptyset$.
- 1110 4. Bulge loop: a special case of interior loop in which there are
 1111 unpaired bases only on one side. i.e $\mathcal{L}_p = \{(i_1, j_1), (i_2, j_2)\}$
 1112 with $i_1 \neq i_2, j_1 \neq j_2$ one of the following assumption holds:
- 1113 • If $\exists i' \in \mathcal{L}_u : i_1 < i' < j_2 \Rightarrow \nexists k' \in \mathcal{L}_u : i_2 < k' < j_2$
- 1114 • If $\exists k' \in \mathcal{L}_u : i_2 < k' < j_2 \Rightarrow \nexists i' \in \mathcal{L}_u : i_1 < i' < j_1$
- 1115 5. Multi-loop: Any loop with degree $d > 2$ i.e. $|\mathcal{L}_p| \geq 3$ and
 1116 $\mathcal{L}_u \neq \emptyset$.
- 1117 6. Exterior loop: a loop in which all the positions are not inter-
 1118 ior of any pair i.e. $\mathcal{L}_p = \emptyset$ and $\mathcal{L}_u \neq \emptyset$.

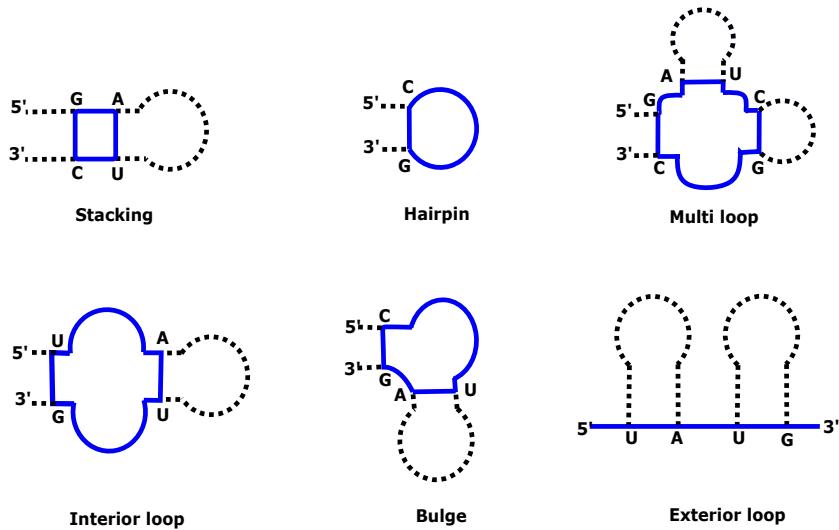


Figure 1.5: RNA secondary structure loop decomposition. Each loop is highlighted in blue.

Definition 5 (Free energy of an RNA secondary structure). Given the loop set $\mathbb{L}_{\phi, \mathcal{S}}$, the free energy ΔG of \mathcal{S} defines its thermodynamic stability. ΔG is the free energy difference with respect to the completely unfolded state [195]. $\Delta G(\mathcal{S}, \phi)$ is computed using

the additivity principle [34], by summing up the energies of its constituent loops.

$$\Delta G(\mathcal{S}, \phi) = \sum_{\mathcal{L} \in \mathbb{L}_{\mathcal{S}, \phi}} \Delta G(\mathcal{L}) \quad (1.2)$$

Many models allow for computing the free energies of those constituent loops, but the dominant is the nearest-neighbor loop energy model [198]. This model associates tabulated free energy values to loop types and nucleotide compositions; the Turner2004 [123] is one of the most widely used parameter sets.

The free energy of each given loop \mathcal{L} is expressed as

$$\Delta G(\mathcal{L}) = \Delta H - T\Delta S \leq 0 \quad (1.3)$$

where ΔH is the (pressure- and volume-dependent) enthalpy change, T the absolute temperature and ΔS the entropy change. The dominant stabilizing effect is attributed to consecutive base-pairs (The stacking loops), whereas long unpaired regions enclosed between base-pairs have destabilizing effects [58, 81]. As a simplified example, the destabilizing free energy contribution $\Delta G(\mathcal{L}_m)$ of a multiloop \mathcal{L}_m as seen in Figure 1.5C is modelled as:

$$\Delta G(\mathcal{L}_m) = \Delta G_{\text{init}} + b\Delta G_{\text{branch}} + u\Delta G_{\text{unpaired}} \quad (1.4)$$

where b is the number of all surrounding base-pairs and u the number of base-pairs [36].

In addition to the definitions mentioned above, we have various properties of an RNA sequence such as structural diversity, positional entropy, structures with maximal expected accuracy, or the density of states. An extensive summary of all possible properties and the history of algorithms is reviewed by Lorenz [115].

The structure decomposition and the tabulated energy parameter sets allow an efficient dynamic programming algorithm to determine a sequence's secondary structure in the entire structure space. Several programs implementing algorithms will enable the computation of these properties efficiently. The thesis gives a literature review of such tools in Chapter 2.

1.5.2 Thermodynamic definitions

A common way to computationally address the RNA folding problem is to consider a dynamic system of structures (the states of

the system). Given enough time, a sequence ϕ will form every possible structure $\mathcal{S} \in \Sigma_\phi$. For each structure $\mathcal{S} \in \Sigma_\phi$, there is a probability of observing it at a given time. This subsection defines RNA folding thermodynamic properties such as structural ensemble, partition function, Boltzmann probability of a structure \mathcal{S} , and the others that derive from them, the base-pair probability and the most probable secondary structure.

The folding tools such as RNAfold, LinearFold used in this thesis use the same thermodynamic definitions. However, some computational folding methods do not rely on a thermodynamic model. For example, Chapter 2 presents a literature review of such tools.

Definition 6 (Structure Ensemble). For a given RNA sequence ϕ , the set of all pseudoknot-free secondary structures with their corresponding energies is called the structure ensemble Σ_ϕ of ϕ or Boltzmann ensemble. We write:

$$\Sigma_\phi = \{\mathcal{S} | \mathcal{S} \text{ is a secondary structure of } \phi\}$$

According to the nearest neighbor energy model, all possible secondary structures of a given RNA sequence do not have the same energy. Since each structure has a unique decomposition, each structure has its own energy but different structures can have the same energy.

Definition 7 (Partition function of RNA). Given the free energy change $\Delta G(\mathcal{S})$ of a structure \mathcal{S} , the partition function $Z(\Sigma_\phi)$ is defined on the Boltzmann ensemble (or structure ensemble) of all possible structures of a given sequence ϕ and we write:

$$Z(\Sigma_\phi) = \sum_{\mathcal{S} \in \Sigma_\phi} \exp(-\beta \Delta G(\mathcal{S}, \phi)) \quad (1.5)$$

Where, $\beta = (RT)^{-1}$ with R the ideal gas constant, and T the temperature.

Definition 8 (Secondary structure probability). How probable is an RNA secondary structure $\mathcal{S} \in \Sigma_\phi$ for the sequence ϕ ? Given the free energy change $\Delta G(\mathcal{S})$ of a structure \mathcal{S} , the boltzmann distribution describes the structure's probability at constant temperature T among all other possible structure of the same sequence ϕ . The probability $p(\mathcal{S} | \phi)$ depends on the free energy $\Delta G(\mathcal{S})$, the lower the more probable. We write:

$$p(\mathcal{S} | \phi) = \frac{\exp(-\beta \Delta G(\mathcal{S}, \phi))}{Z} \quad (1.6)$$

1175 where, Z is the partition function and $\beta = (RT)^{-1}$ the thermal
 1176 constant.

1177 **Definition 9** ([MFE](#) secondary structure). To predict biologically
 1178 relevant structures, most computational methods search for struc-
 1179 tures that minimize the free energy. For a given sequence ϕ , let Σ_ϕ
 1180 be the secondary structure ensemble of ϕ . The minimum free en-
 1181 ergy structure \mathcal{S}_{MFE} is the structure with the lowest probability
 1182 $p(\mathcal{S}|\phi)$ i.e. the most stable conformation in the thermodynamic
 1183 equilibrium. We write:

$$\mathcal{S}^{MFE}(\phi) = \arg \min_{\mathcal{S} \in \Sigma_\phi} \Delta G(\mathcal{S}, \phi) \quad (1.7)$$

1184 **Definition 10** (Base-pair probability). Let $\phi = (\phi_i)_{1 \leq i \leq L}$ be an
 1185 [RNA](#) sequence. The base-pair probability matrix $\mathbf{P}(\phi)$ quantifies
 1186 the equilibrium structural features of the ensemble Σ_ϕ , with en-
 1187 tries $P_{i,j}(\phi) \in [0, 1]$ defines as follows:

$$P_{i,j}(\phi) = \sum_{\mathcal{S} \in \Sigma_\phi} p(\mathcal{S}|\phi) S_{i,j}(\mathcal{S}) \quad (1.8)$$

1188 $P_{i,j}(\phi)$ corresponds to the probability that base-pair i,j forms
 1189 at the equilibrium. $\mathbf{S}(\mathcal{S})$ is the structure matrix with entries $S_{i,j} \in$
 1190 $\{0, 1\}$. If the structure \mathcal{S} contains pair (i, j) , then $S_{i,j}(\mathcal{S}) = 1$ other-
 1191 wise $S_{i,j}(\mathcal{S}) = 0$.

1192 The base pair probabilities enable then a new view at the struc-
 1193 ture ensemble. [Figure 1.6](#) shows an example of [MFE](#) structure and
 1194 the base pair probability dot plot¹ of a [tRNA](#). A square at row i
 1195 and column j indicates a base pair. The area of a square in the
 1196 upper right half of the matrix is proportional to the base pair
 1197 probability (i, j) within the equilibrium ensemble. The lower left
 1198 half shows all pairs belonging to the [MFE](#) structure. While the [MFE](#)
 1199 consists of hairpins, bulge and stacking, several different loops
 1200 are visualized in the pair probabilities, which leads to several
 1201 local minima with different shapes.

¹ computed using [RNAfold -p](#)

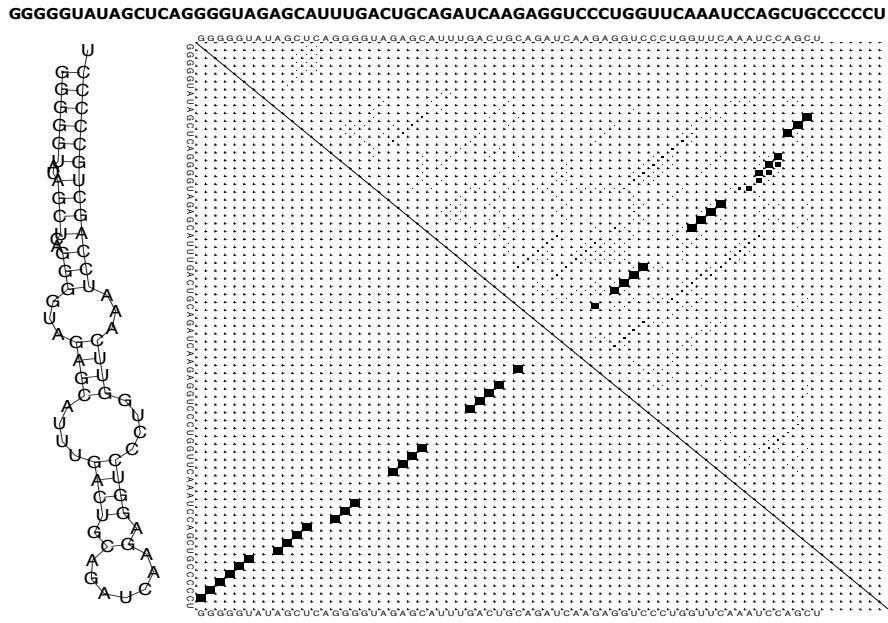


Figure 1.6: **Base-pair probability matrix** of a tRNA sequence computed using RNAfold 2.4.13. The MFE structure is depicted on the left and the sequence on top. The frequency of the MFE structure in the structural ensemble Σ_ϕ is 0.116. The dot plot on the right shows the pair probabilities within the equilibrium ensemble as (72×72) -matrix and is an excellent way to visualize structural alternatives.

1202 The definitions mentioned above provide us with a necessary
 1203 framework enabling us to compute the MFE secondary structure
 1204 within the equilibrium ensemble Σ_ϕ . Several implementations
 1205 of these definitions have been suggested [112, 148, 218], and
 1206 they are available as an application programming interface (API).
 1207 In the context of this work, we are not only interested in the
 1208 MFE structure but, instead, we use some features of the existing
 1209 computer libraries (e.g. the computation of the structure free
 1210 energy) to predict an ensemble structure. The following section
 1211 introduces some metrics used in this dissertation to compare RNA
 1212 secondary structures and, eventually, the structure predictions
 1213 produced by different tools.

1214 1.5.3 *Structural distance definitions*

1215 The validation of the results obtained in this thesis is purely
 1216 empirical. We achieved this goal by comparing the predicted and
 1217 expected structures for the folding tools. We use the PPV and
 1218 the sensitivity's statistical properties for the benchmark results

1219 presented in [Chapter 3](#). For the inverse folding tools, we compare
 1220 the [MFE](#) structure of the designed sequence to the target structure.
 1221 For that end, a rigorous definition of a measure of similarities
 1222 between two structures is needed. This subsection defines the
 1223 different similarity measurements used throughout this work. In
 1224 addition, it defines the objective functions used in our inverse
 1225 folding presented in [Chapter 5](#).

1226 **Definition 11** (The [PPV](#)). It measures the fraction of correct base-
 1227 pairs in the predicted structure and it is defined as follows:

$$\text{PPV} = \frac{TP}{TP + FP} \quad (1.9)$$

1228 where TP and FP stand respectively for the number of correctly
 1229 predicted base-pairs (true positives), and the number of wrongly
 1230 predicted base-pairs (false positives).

Definition 12 (Sensitivity). It measures the fraction of base-pairs
 in the accepted structure that are predicted.

$$\text{Sensitivity} = \frac{TP}{TP + FN} \quad (1.10)$$

1231 where FN stands for the number of base-pairs not detected (false
 1232 negatives).

Definition 13 (Base-pair distance). Let σ_1 and σ_2 be two sec-
 ondary structures in their string representation. The base-pair
 distance between σ_1 and σ_2 is defined as follows:

$$d_{bp}(\sigma_1, \sigma_2) = \sum_{i,j} A_{i,j}[\sigma_1] + A_{i,j}[\sigma_2] + 2 \times A_{i,j}[\sigma_1]A_{i,j}[\sigma_2], \quad (1.11)$$

where,

$$A_{i,j}[\sigma] = \begin{cases} 1 & \text{if } (i,j) \text{ is a base-pair in } \sigma \\ 0 & \text{otherwise} \end{cases}$$

1233 **Definition 14** (Hamming distance). Let σ_1 and σ_2 be two sec-
 1234 ondary structures in their string representation. We define the
 1235 hamming distance between σ_1 and σ_2 , $d_h(\sigma_1, \sigma_2)$, to be the num-
 1236 ber of position where σ_1 and σ_2 differ.

$$d_h(\sigma_1, \sigma_2) = \sum_{i=1}^L S(\sigma_1^i, \sigma_2^i) \quad (1.12)$$

where,

$$S(\sigma_1^i, \sigma_2^j) = \begin{cases} 1 & \text{if } \sigma_1^i \neq \sigma_2^j \\ 0 & \text{otherwise} \end{cases}$$

Definition 15 (Ensemble defect (ED)) [219]). Given an RNA sequence ϕ of length L , the ensemble defect \mathcal{D}_E is the expected base-pair distance between a target structure \mathcal{S}^* and a random structure generated with respect to the Boltzmann probability distribution. It is defined as follows:

$$\begin{aligned} \mathcal{D}_E(\phi, \mathcal{S}^*) &= \sum_{\mathcal{S} \in \Sigma_\phi} p(\mathcal{S}|\phi) d_{bp}(\mathcal{S}, \mathcal{S}^*) \\ &= L - \sum_{1 < i, j < L} P_{i,j}(\phi) S_{i,j}(\mathcal{S}^*) \end{aligned} \quad (1.13)$$

1237 where $P_{i,j}$ is the base-pair probability matrix entrances, $d_{bp}((\mathcal{S}, \mathcal{S}^*))$
1238 is the base pair distance between two structures, and $\mathbf{S}(\mathcal{S})$ is the
1239 structure matrix with entries $S_{i,j} \in \{0, 1\}$. If the structure \mathcal{S} con-
1240 tains pair (i, j) , then $S_{i,j}(\mathcal{S}) = 1$ otherwise $S_{i,j}(\mathcal{S}) = 0$.

1241 **Definition 16** (Normalized Energy Distance (NED)). the differ-
1242 ence between the energy of a given sequence ϕ evaluated to fold
1243 into a target structure \mathcal{S}^* and the minimum free energy of the
1244 sequence in its structural ensemble Σ_ϕ . The value is normalized
1245 over all the sequences in a given population P .

$$\mathcal{N}_E(\phi, \mathcal{S}^*) = [1 - \Delta \hat{E}(\mathcal{S}^*, \phi)]^q \quad \forall q > 1 \quad (1.14)$$

where,

$$\Delta \hat{E}(\mathcal{S}^*, \phi) = \frac{\Delta E(\mathcal{S}^*, \phi)}{\sum_{s \in P} \Delta E(\mathcal{S}^*, s)} \quad (1.15)$$

and,

$$\Delta E(\mathcal{S}^*, \phi) = \Delta G(\mathcal{S}^*, \phi) - \arg \min_{\mathcal{S} \in \Sigma_\phi} \Delta G(\mathcal{S}, \phi) \quad (1.16)$$

1246 Among the definitions mentioned above, 11 and 12 are used in
1247 Chapter 3 for the benchmark comparison. Whereas, definitions
1248 13, 14, 15, 16 are used in Chapter 5 for both objective function and
1249 benchmark purposes. The following section provides a formal
1250 definition of the fitness landscape and some of its properties. It
1251 will mostly use 14 for both structure and sequence comparison.

1252 1.5.4 RNA folding map properties

1253 This work considers RNA molecule folding and inverse folding
 1254 optimisation problems. In both cases, It is fundamental to define
 1255 the fitness landscape notion. This subsection provides the formal
 1256 definitions of the fitness landscape and examples related to the
 1257 folding and inverse problem. Some properties such as neutrality,
 1258 mutation mode or move operator are also provided. The size
 1259 of the RNA structural ensemble has been analytically computed
 1260 through tools developed by Stein and Waterman [185], and it
 1261 yields an upper bound of $S_L \approx 1.48 \times L^{-\frac{3}{2}} 1.85^L$ structure vis-a-
 1262 vis 4^L sequences. Compared to the total number of sequences,
 1263 the number of structures is much smaller, which means there is
 1264 a high possibility that many sequences fold into the same MFE
 1265 secondary structure. In case that happens, we call the set of those
 1266 sequences a neutral set. The fraction of such sequences defines
 1267 the neutrality of a fitness landscape.

1268 **Definition 17** (Fitness landscape). A fitness landscape \mathcal{L} results
 1269 from the combination of three elements: a set of configurations
 1270 \mathcal{V} , a cost or fitness function f , and a *move* operator ψ that induces
 1271 a topology on the set of configurations. We write:

$$\mathcal{L} = (\mathcal{G}_f, f, \psi) \quad (1.17)$$

1272 where \mathcal{G}_f is the the landscape underlying the hypergraph whose
 1273 vertices are the elements from \mathcal{V} labelled with values given by f ,
 1274 and whose edges are specified by the move operator ψ .

1275 The fitness function f assigns to each configuration $v \in \mathcal{V}$ a real
 1276 value taken from an interval $\mathbb{I} \subset \mathbb{R}$ as follows:

$$f : \mathcal{V} \rightarrow \mathbb{I}$$

1277 An example of fitness function in the case of inverse folding
 1278 is defined in Chapter 5 (Section 5.1.2), which uses the hamming
 1279 distance d_h and $\mathcal{V} = \{A, C, G, U\}^L$. But in this case, the fitness
 1280 defined in the structural space Σ_ϕ . i.e. we have an intermediate
 1281 folding function $\Delta G(\mathcal{S}, \phi)$, mapping any sequence $\phi \in \mathcal{V}$ to an
 1282 MFE secondary structure.

The move (or mutation) operator ψ defines the relationship
 between the configuration from \mathcal{V} in the following way:

$$\psi : \mathcal{V} \rightarrow \mathcal{V}$$

1283 **Definition 18** (Mutation mode). Let $\phi, \phi' \in \mathcal{V} = \{A, C, G, U\}^L$,
 1284 be two RNA sequences. ϕ' is said to be an n -point mutation of ϕ if
 1285 it differs from ϕ at n nucleotides; i.e. $d_h(\phi, \phi') = n$ where $d_h(., .)$
 1286 is the hamming distance on $\{A, C, G, U\}^L$.

1287 A mutation mode is a random variable U taking values in
 1288 $\{1, \dots, L\}$. $P(U = n)$ is defined as the probability that, exactly n
 1289 nucleotides, selected uniformly at random undergo point muta-
 1290 tion during a mutation event. U can generally be any probability
 1291 distribution.

1292 **Definition 19** (Neutral set of RNA sequences). For a give fitness
 1293 landscape $\mathcal{L} = (\mathcal{G}_f, f, \psi)$, with $\mathcal{V} = \{A, C, G, U\}^L$, two RNA se-
 1294 quence ϕ_1 and ϕ_2 are set to be neutral $\iff f(\phi_1) = f(\phi_2)$. We
 1295 call a set $\Gamma \subset \mathcal{V}$ of all such RNA sequences a neutral set. In the
 1296 case of inverse folding, ϕ_1 and ϕ_2 are neutral if they share the
 1297 same MFE secondary structure. In contrast, ϕ_1 and ϕ_2 have the
 1298 same free energy in the folding problem context.

1299 **Definition 20** (Neutral Network). Let $\mathcal{G}(\mathcal{V}, E)$ be a connected
 1300 graph in which vertices are all in the neutral sequence set Γ (i.e.
 1301 $\mathcal{V} \subset \Gamma$). \mathcal{G} is said to be a neutral network $\iff \forall e(v_i, v_j) \in E, v_i, v_j$
 1302 differ by a single nucleotide (i.e. $d_h(v_i, v_j) = 1$).

1303 We provided in this subsection a general definition of a fitness
 1304 landscape with examples related to computational RNA folding
 1305 and inverse folding. Now that we have all the ingredients to
 1306 computationally study the folding and the inverse folding of RNA
 1307 molecule, we are left with the definition of some computational
 1308 techniques used in our proposed tools. Our contributions rely on
 1309 two well-known techniques of algorithms: the FFT for the folding
 1310 mechanism and the EA for the inverse folding. An overview of
 1311 both techniques is provided in the following section.

1312 1.5.5 *The fast Fourier transform (FFT) and evolutionary algorithm 1313 (EA) applied to RNA bioinformatics*

1314 The computational results present in this work rely on two well-
 1315 known techniques: the FFT and EA. Both approaches have already
 1316 been studied and have found many applications, including the
 1317 computational folding and inverse folding of ncRNA. This section
 1318 gives a short overview of the two concepts.

1319 A FFT is an algorithm that computes the distcrete Fourier trans-
 1320 form (DFT) of a sequence or its inverse (Inverse discrete Fourier
 1321 transform (IDFT)). Fourier analysis converts a signal from its

original domain (often time or space) to a representation in the frequency domain and vice versa. The **DFT** is obtained by decomposing a sequence of values into components of different frequencies.

More formally, let $\{x_k\} := x_0, \dots, x_{L-1}$ be a sequence of L complex numbers, the **DFT** transforms the sequence $\{x_k\}$ into another sequence of L complex numbers $\{X_k\} := X_0, \dots, X_{L-1}$ defined as follows:

$$X_k = \sum_{n=0}^{L-1} x_n e^{-i2\pi kn/N} \quad (1.18)$$

The direct evaluation of [Equation 1.18](#) will require $O(L^2)$ operations because there are L outputs of X_k , and each of them requires a sum of L terms. A **FFT** is, therefore, any approach allowing to compute the same results in $O(L \log L)$ operations [92].

Let x and y be two sequences of length L and let X and Y be their respective **DFTs**. The correlation c_k between sequence x and y with the positional lag of k sites is defined as follows:

$$c_k = \sum_{1 < n < L, 1 < n+k < L} x_n y_{n+k} \quad (1.19)$$

It is known that the correlation c_k can be expressed in terms of the **DFT**. We write:

$$c_k \Leftrightarrow X_n^* \cdot Y_n \quad (1.20)$$

Where the asterisk denotes complex conjugation. That means we simply need to compute the **DFT** X_n and Y_n . Therefore, we can compute correlations c_k using the **FFT** as follows: **FFT** the two sequences, multiply one resulting transform by the complex conjugate of the other, and inverse transform the product.

Similar to [Equation 1.18](#), the direct evaluation of c_k requires $O(L^2)$ operations and taking advantage of the **FFT** reduces it to $O(L \log L)$ operations. Several **FFT** algorithms have been implemented to speed up the computation of the **DFT** but so far, the most commonly used is the Cooley–Tukey algorithm [25].

The same idea has been applied in the context of **RNA** bioinformatics, where the two sequences of complex numbers can be thought of as two data sets of real numbers encoding the **RNA** sequences information. And the correlation c_k measures the homologous region in the two **RNA** sequences [94]. In contrast to

1352 Katoh and his collaborators [94], we use the **FFT** to rapidly iden-
 1353 tify the largest stems of an **RNA** sequence. Thanks to the **FFT** which
 1354 allows us to efficiently predict the fast-folding pathways of **RNA**
 1355 molecules (See [Chapter 3](#)) within a reasonable **CPU** time.

1356 The **EA** is another well-known heuristic approach, especially
 1357 when dealing with problems in which less information about
 1358 the fitness landscape is provided or when there is no exact al-
 1359 gorithm in polynomial for such problems. The **EA** approach is
 1360 inspired by evolutionary systems. In the 1950s and the 1960s,
 1361 several computer scientists already independently studied evo-
 1362 lutionary systems with the idea that evolution could be used as
 1363 an optimization tool for engineering problems [131]. The picture
 1364 in all these systems was to evolve a population of candidate so-
 1365 lutions to a given situation, using operators inspired by natural
 1366 genetic variation and natural selection.

1367 Since the genetic algorithm (or more generally **EA**) was pro-
 1368 posed by John Holland [82] in the early 1970s, it has emerged
 1369 as a popular search heuristic. It has found application in many
 1370 disciplines that deal with complex landscape optimization prob-
 1371 lems, e.g. RNA folding [132, 211] and inverse RNA folding [47,
 1372 48, 189].

1373 **EAs** form a class of heuristic search methods based on a par-
 1374 ticular algorithmic framework whose main components are the
 1375 variation operators (mutation and recombination or crossover)
 1376 and the selection operators (parent selection and survivor selec-
 1377 tion). The general evolutionary algorithm framework is depicted
 1378 in [Figure 1.7](#). In most of the **EA** implementations, the solutions
 1379 are encoded in the form of genomes (array of elements). The
 1380 simplest form of **EA** typically involves two types of operators:
 1381 selection and mutation (single point).

- 1382 • **selection:** the operator consists of selecting solutions in
 1383 the population for reproduction. The fitter the solution,
 1384 the more times as likely it is selected to reproduce. This
 1385 operator often requires a fitness function evaluation.
- 1386 • **mutation:** the operator allows generating new solutions in
 1387 the population. It randomly flips or permutes some element
 1388 positions in a genome solution. For example, if we encode
 1389 the solutions in a binary string, the solution 00000100 might
 1390 be mutated in its second position to yield 01000100. Mu-
 1391 tation can occur at each bit position in a string with some
 1392 probability, usually very small (e.g. 0.001 for a sequence of
 1393 length 50).

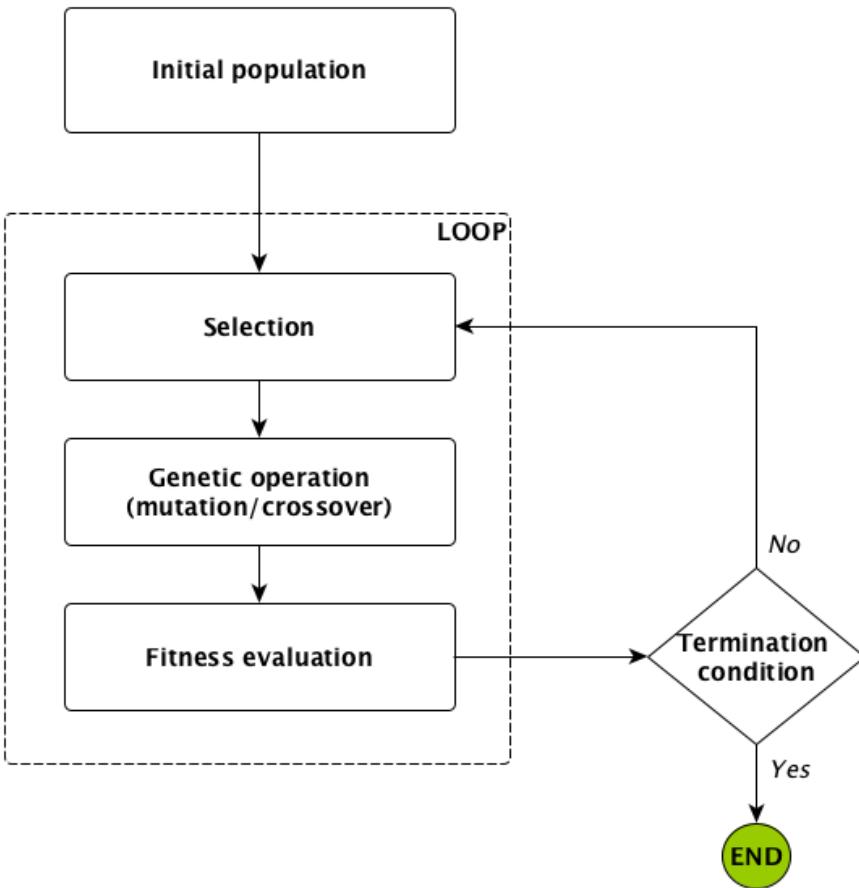


Figure 1.7: **Evolutionary algorithm flow diagram.** The algorithm initializes a population of candidate solutions and then loops over the three genetic operations until the termination criteria are satisfied.

1394 In a more complex configuration, we can have a crossover operator
 1395 that plays almost the same role as mutation, which generates
 1396 new solutions in the population. In contrast to the mutation operator,
 1397 the crossover randomly chooses a locus and exchanges the
 1398 subsequence solutions before and after that locus between two
 1399 solutions to create two offspring solutions.

1400 In the context of this work, we use the simplest form of EA, in
 1401 which we did not consider a crossover operator. We implement
 1402 the simplest EA framework with a mutation operator adapted
 1403 to the Inverse folding problem, which results in an alternative
 1404 computational tool named aRNAque (see Chapter 5)).

1405 This section provided an overview of the two main tools used
 1406 in this thesis, which are EA and FFT. The EA is implemented in
 1407 the computational inverse folding tool we propose in Chapter 5,

¹⁴⁰⁸ and the [FFT](#) in the [RNA](#) folding tool that will be introduced in
¹⁴⁰⁹ [Chapter 3](#).

¹⁴¹⁰ 1.6 CONCLUSION AND OUTLINE OF THE THESIS

¹⁴¹¹ This introductory chapter presents nucleic acids in general and,
¹⁴¹² in particular, a description of [ncRNA](#) and its chemical, biologi-
¹⁴¹³ cal, and algorithmic definitions. Those concepts with biological
¹⁴¹⁴ motivations constitute the basis of the thesis.

¹⁴¹⁵ We organize the next part of the thesis into fives. The two first
¹⁴¹⁶ s are grouped into a first result part which only concerns [RNA](#)
¹⁴¹⁷ folding. The second part discusses inverse folding, and similarly
¹⁴¹⁸ to the first part, it contains two s. The last discusses the presented
¹⁴¹⁹ results and concludes by providing some limitations and possible
¹⁴²⁰ future research directions.

¹⁴²¹ In [Part i](#), [Chapter 2](#) provides a brief literature review on the
¹⁴²² existing computational methods for [RNA](#) folding. The review fo-
¹⁴²³ cuses on thermodynamic and machine learning methods such
¹⁴²⁴ as [RNAfold](#), [LinearFold](#) and [Mfold](#). We review some of the limi-
¹⁴²⁵ tations of existing tools in [Chapter 2](#), such as the computational
¹⁴²⁶ time, and in some cases, the predicted thermodynamic structure
¹⁴²⁷ does not match the native one. [Chapter 3](#) presents our proposed
¹⁴²⁸ folding tool called [RAFFT](#), which aims at overcoming those lim-
¹⁴²⁹ itations. [RAFFT](#) implements a novel heuristic to predict [RNA](#) sec-
¹⁴³⁰ ondary structure formation pathways that has two components:
¹⁴³¹ (i) a folding algorithm and (ii) a kinetic ansatz. This heuristic is in-
¹⁴³² spired by the kinetic partitioning mechanism, by which molecules
¹⁴³³ follow alternative folding pathways to their native structure, some
¹⁴³⁴ much faster than others. [RAFFT](#) starts by generating an ensemble
¹⁴³⁵ of concurrent folding pathways ending in multiple metastable
¹⁴³⁶ structures, which contrasts with traditional thermodynamic ap-
¹⁴³⁷ proaches that find single structures with minimal free energies.
¹⁴³⁸ When analyzing 50 predicted folds per sequence, we found near-
¹⁴³⁹ native predictions for [RNAs](#) of length ≤ 200 nucleotides, match-
¹⁴⁴⁰ ing the performance of current deep-learning-based structure
¹⁴⁴¹ prediction methods [[160](#), [220](#)]. [RAFFT](#) also acts as a folding ki-
¹⁴⁴² netic ansatz, which we tested on two [RNAs](#): the [CFSE](#) and a classic
¹⁴⁴³ bi-stable sequence. For the [CFSE](#), an ensemble of 68 distinct struc-
¹⁴⁴⁴ tures computed by [RAFFT](#) allowed us to produce complete folding
¹⁴⁴⁵ kinetic trajectories. In contrast, known methods require evaluat-
¹⁴⁴⁶ ing millions of sub-optimal structures to achieve this result. For
¹⁴⁴⁷ the second application, only 46 distinct structures were required

¹⁴⁴⁸ to reproduce the kinetics, whereas known methods required a
¹⁴⁴⁹ sample of 20,000 structures.

¹⁴⁵⁰ Similar to the first part of the result, [Part ii](#) contains two chapters.
¹⁴⁵¹ [Chapter 4](#) will briefly introduce the [RNA](#) design problem. It
¹⁴⁵² distinguishes the positive from the negative [RNA](#) design problem
¹⁴⁵³ and reviews the current state of the art computational tools, espe-
¹⁴⁵⁴ cially those implementing evolutionary techniques. The existing
¹⁴⁵⁵ tools present challenges when benchmarked on recent datasets
¹⁴⁵⁶ such as [Eterna100](#). Another limitation is that most existing tools
¹⁴⁵⁷ do not consider the pseudoknot patterns in their designing pro-
¹⁴⁵⁸ cess. In [Chapter 5](#), we propose an improved evolutionary algo-
¹⁴⁵⁹ rithm inspired by the Lévy flights. Like a Lévy flight, our tool,
¹⁴⁶⁰ [aRNAque](#), implements a Lévy mutation scheme that allows simul-
¹⁴⁶¹ taneous search at all scales over the mutational landscape. New
¹⁴⁶² mutations often produce nearby sequences (one-point mutations)
¹⁴⁶³ but occasionally generate mutant sequences far away in geno-
¹⁴⁶⁴ type space (macro-mutations). In [aRNAque](#), the number of point
¹⁴⁶⁵ mutations distribution at every step is taken to follow a Zipf
¹⁴⁶⁶ distribution. The Lévy mutation scheme increases the diversity
¹⁴⁶⁷ of designed [RNA](#) sequences and reduces the average number of
¹⁴⁶⁸ evaluations of the evolutionary algorithm compared to the local
¹⁴⁶⁹ search. The overall performance showed improved empirical re-
¹⁴⁷⁰ sults compared to existing tools through intensive benchmarks on
¹⁴⁷¹ both pseudoknot (the [PseudoBase++](#) dataset) and pseudoknot-
¹⁴⁷² free (the [Eterna100](#) dataset) datasets.

¹⁴⁷³ Finally, [Chapter 7](#) presents a general conclusion, a discussion
¹⁴⁷⁴ on the results obtained and some promising perspectives. It em-
¹⁴⁷⁵ phasizes the understanding of the Lévy mutation in the context
¹⁴⁷⁶ of [RNA](#) design and the application of our results to evolutionary
¹⁴⁷⁷ dynamics.

[August 4, 2022 at 19:02 – 1.0]

1478

Part I

1479

RNA FOLDING

1480

This first part of our thesis provides a literature review on existing computational tools addressing the prediction of RNA secondary structure, and it presents our proposed tool RAFFT. Chapter 3 contains figures and ideas that have previously appeared in our publication:

1486

- [138] Vaitea Opuu, Nono SC Merleau, Vincent Messow and Matteo Smerlak (2021). RAFFT: Efficient prediction of RNA folding pathways using the fast Fourier transform. In: *bioRxiv* (Submitted and accepted) (PLoS Comp. Biol.)

1487

1488

1489

1490

[August 4, 2022 at 19:02 – 1.0]

2

1491

1492 INTRODUCTION TO RNA FOLDING

1493 We provided some motivations for studying ncRNAs and intro-
1494 duced their bioinformatic concepts in the introduction. We also
1495 highlighted the relationship between the structure of ncRNAs and
1496 their functions. The functions of ncRNAs and their lengths usu-
1497 ally distinguish them, and several ncRNA classes were presented.
1498 Identifying the ncRNA functions is challenging, though there is
1499 a widespread expectation that their functions are largely deter-
1500 mined by their structures. The process of determining the RNA
1501 structure is often termed RNA folding. Experimental methods
1502 that determine the secondary structure of such molecules are
1503 usually expensive. Many computational methods have been de-
1504 veloped in the last decades as alternatives. This chapter overviews
1505 computational methods for predicting RNA secondary structures.
1506 Two techniques will be reviewed: statistical approaches such as
1507 machine learning and score-based methods.

1508 **2.1 STABILITY AND PREDICTION OF RNA SECONDARY STRUC-**
1509 **TURES**

1510 The mapping from RNA sequences to their corresponding sec-
1511 ondary structure defines the folding of RNA molecules. RNA fold-
1512 ing is, therefore, a process by which a linear RNA sequence ac-
1513quires a secondary structure through intra-molecular interactions.
1514 The nature of those interactions defines the thermodynamic sta-
1515 bility of the secondary structure. Throughout this dissertation, we
1516 will denote the thermodynamic stability of a structure σ by ΔG_σ ,
1517 which is the free energy difference with respect to the completely
1518 unfolded state. This section provides an intuition on how the *free*
1519 *energy* of an RNA secondary structure is computed based on the
1520 definitions and concepts introduced in Chapter 1. Furthermore,
1521 it introduces the problem of RNA secondary structure prediction
1522 and an overview of existing techniques.

1523 In predicting biologically relevant structures, most computa-
1524 tional methods search for structures that minimize the free energy
1525 function ΔG (i.e. the MFE structure). Therefore, the prerequisite
1526 to efficiently computing the MFE secondary structure is the com-
1527 putation of the free energy for any given secondary structure

1528 \mathcal{S} . The calculation of the [RNA](#) structure free energies starts by
 1529 decomposing each structure into components called loops (See
 1530 Definition 4). The loop decomposition allows building the ba-
 1531 sis of the standard energy model for [RNA](#) secondary structures
 1532 called the nearest neighbour ([NN](#)) model [198]. The total free
 1533 energy of a secondary structure is assumed to be a sum over its
 1534 constituent loops according to the additivity principle [34] (see
 1535 Definition 5). Therefore, this structure decomposition allows an
 1536 efficient dynamic programming ([DP](#)) algorithm to determine the
 1537 [MFE](#) pseudoknot-free structure of a sequence ϕ in the structure
 1538 space Σ_ϕ .

1539 The [DP](#) is a computer programming method developed by
 1540 Richard Bellman in the early 1950s [10], and it has found appli-
 1541 cations in various fields, including the [RNA](#) secondary structure
 1542 prediction. It consists of simplifying a complicated problem by
 1543 breaking it down into simpler sub-problems in a recursive man-
 1544 ner. When sub-problems can be nested recursively inside larger
 1545 problems so that [DP](#) methods are applicable, then there is a re-
 1546 lation between the value of the larger problem instance and the
 1547 values of the sub-problems.

1548 For example, let us consider the definition of secondary struc-
 1549 ture \mathcal{S} introduced in the previous chapter (Definition 2) and its
 1550 string representation σ . When considering a substructure $\sigma[i:j]$
 1551 within the sequence interval $\phi[i:j]$, there are only two alterna-
 1552 tives to how position i may contribute to $\sigma[i:j]$. Either i does not
 1553 pair with any other position, or it pairs with another nucleotide k
 1554 with $i < k \leq j$. In the first situation, $\sigma[i:j]$ consists of the base
 1555 pairs in the subsequence $\sigma[i+1:j]$ only. The formation of a
 1556 base pair (i,k) , however, subdivides the structure into two parts,
 1557 one enclosed by (i,k) , namely $\sigma[i+1:k-1]$, and the other one,
 1558 $\sigma[k+1:j]$. Thus, $\mathcal{S} = \text{proc}\{\sigma[i+1:k-1] \cup \sigma[k+1:j]\} \cup \{(i,k)\}$,
 1559 where the *proc* is the recursive procedure. Since condition (3) of
 1560 definition 2 ensures that the position (i,j) can not contain base
 1561 pairs that cross (i,k) (or at least in the pseudoknot-free situation),
 1562 the two shorter substructures $\sigma[i+1:k-1]$ and $\sigma[k+1:j]$ can
 1563 be treated independently for a large variety of purposes.

1564 This observation has led to a recursive decomposition scheme
 1565 for [RNA](#) secondary structures, which is the basis of the large
 1566 variety of [DP](#) approaches that solve [RNA](#) secondary structure pre-
 1567 diction problems. The first [DP](#) algorithm was then proposed by
 1568 Nussinov and Jacobson [136] to find the structure with the max-
 1569 imum base pairs. A few years after, Zucker and Stieger [228]
 1570 extended Nussinov's algorithm to a more realistic scoring model

1571 based on free energy, the [NN](#) model. Almost all score-based meth-
1572 ods rely on the same [DP](#) algorithm, but the decomposition scheme
1573 and the scoring model could differ from one to another. When
1574 predicting structures with non-canonical base pairs, some other
1575 scoring schemes are used, such as nucleotide cyclic motifs score
1576 system [[28](#), [140](#), [174](#)] or equilibrium partition function [[177](#)].

1577 In addition to score-based methods, we have comparative se-
1578 quence analysis methods, which are the most computationally
1579 accurate for determining [RNA](#) secondary structures [[72](#), [119](#)]. Us-
1580 ing the set of homologous structures, the comparative method
1581 allows finding base pairs that covary to maintain [WC](#) and wobble
1582 bases of a given sequence ϕ [[73](#)]. The first comparative method
1583 predicting a common secondary structure conserved in the given
1584 homologous sequence set was developed by Han and Kim in
1585 the early nineteenth century, and it was based on comparative
1586 phylogenetic analysis.

1587 When neglecting the special base pairs (or pseudoknots) and
1588 the weak interactions, the running time of both approaches (score-
1589 based and comparative analysis) is usually $O(L^3)$ (Where L is
1590 the [RNA](#) sequence length) and thus prohibitingly slow for longer
1591 sequences.. Many other comparative analysis methods and vari-
1592 ations of score-based methods were also proposed to improve
1593 computational time. More recently, a heuristic method such as
1594 [LinearFold](#) allows achieving good [RNA](#) folding performance in
1595 a linear time ($O(L)$).

1596 When pseudoknots are considered, the loop decomposition
1597 of a secondary structure and the energy rules break down. Al-
1598 though we can assign reasonable free energies to the helices in a
1599 pseudoknot and even to possible coaxial stacking between them,
1600 it is impossible to estimate the effects of the new kinds of loops
1601 created. Base triples pose an even greater challenge because the
1602 exact nature of the triple cannot be predicted in advance, and even
1603 if it could, we have no data for assigning free energies. Neverthe-
1604 less, there are existing techniques that approximate the energies
1605 of pseudoknot loops and allow the dynamic programming tech-
1606 nique to tackle the [RNA](#) folding with pseudoknots. However, the
1607 time complexity still remains the main problem. Using the [DP](#)
1608 technique for the pseudoknot structure prediction, the time com-
1609 plexity goes up to $O(L^6)$ for the exact prediction. But for heuristic
1610 methods such as [IPKnot](#) [[162](#)] and [Hotknots](#) [[147](#)], the running
1611 time can be reduced down to $O(L^4)$.

1612 Despite the advanced development of computational tools for
1613 [RNA](#) folding, it's challenging to understand the folding mecha-

1614 nism fully. In contrast to score-based and comparative analysis
1615 methods, machine learning methods are data-driven methods
1616 that require no knowledge of the folding mechanism. Neverthe-
1617 less, the requirement of **ML**-based methods is a large amount of
1618 training data on which they can learn. In the last few decades,
1619 **ML** methods have been used for many aspects of **RNA** secondary
1620 structure prediction methods to improve the prediction perfor-
1621 mance and overcome the limitations of existing methods. How-
1622 ever, they did not replace the mainstream score-based methods
1623 with respect to accuracy and generalization. In addition to some
1624 overfitting concerns, **ML**-based methods cannot give dynamic
1625 information on the **RNA** folding process since little data are avail-
1626 able on structural dynamics. In addition, the training data used
1627 in **ML**-based methods are mostly obtained through phylogenetic
1628 analyses. Consequently, their prediction may be biased due to the
1629 *in vivo* third elements. The following subsections provide a de-
1630 tailed description of some of the recent **ML**-based and score-based
1631 tools for secondary structure prediction.

1632 In sum, computational methods usually consider the **MFE** sec-
1633 ondary structure as the most biologically relevant one. Predicting
1634 the **MFE** structure consists of solving a free energy optimization
1635 problem in the case the scoring function is the free energy. Ex-
1636 isting methods for RNA secondary structures prediction can be
1637 clustered into three main categories: the scored-based, compara-
1638 tive sequence analysis and **ML** methods. The **DP** technique is
1639 one of the most widely used score-based methods, but they are
1640 usually less accurate than the comparative methods. In contrast,
1641 **ML** methods are more recent and still under intensive improve-
1642 ments. The following section will overview some existing tools
1643 and highlight their limitations.

1644 **2.1.1 MFE prediction tools for pseudoknot-free **RNA** sequences using**
1645 **a score-base method**

1646 The score-based methods often assume that the native or bio-
1647 logical **RNA** structure is the one that minimizes/maximizes the
1648 overall total score, depending on the hypotheses made on the
1649 **RNA** folding mechanism. In the pseudoknot-free **MFE** prediction,
1650 where the special and weak interactions are neglected, the fold-
1651 ing problem is less complex, and the scoring model is simply the
1652 free energy. Hence, the issue of **RNA** secondary structure predic-
1653 tion becomes an optimization problem that aims at finding the

1654 best-scoring structure \mathcal{S}^{MFE} by minimizing a scoring function
 1655 ΔG .

$$\mathcal{S}^{MFE} = \operatorname{argmin}_{\mathcal{S} \in \Sigma_{\phi}} \Delta G(\mathcal{S}, \phi) \quad (2.1)$$

1656 Where Σ_{ϕ} is the set of all possible pseudo-knot free secondary
 1657 structures for the sequence ϕ of length L and, $\Delta G(\mathcal{S}, \phi)$ the free
 1658 energy of the structure \mathcal{S} evaluated for the sequence ϕ .

1659 Since each possible structure can be uniquely and recursively
 1660 decomposed into smaller components (or loops) with indepen-
 1661 dent free energy contributions, the DP is best suited for most of
 1662 the following tools presented here.

- 1663 • **Unfold** [227, 228]: It is the successor of the original mfold
 1664 program which was the first realistic implementation of the
 1665 DP for secondary structure predictions with a score based
 1666 on the loop energy parameters and a worse case time com-
 1667 plexity of $O(L^3)$. The initial version was an improvement of
 1668 the simplest DP for secondary structure prediction known as
 1669 the *maximum circular matching problem*. The authors demon-
 1670 strated that the loop-based energy model is also amenable
 1671 to the same algorithmic ideas. With McCaskill's algorithm
 1672 [127], for computing the partition function of the equilib-
 1673 rium ensemble of RNA molecules, more efficient implemen-
 1674 tations of the initial program with accurate thermodynamic
 1675 modelling have been provided. The latest implementation
 1676 is known as Unfold.
- 1677 • **RNAstructure** [125, 148]: The software first appeared in
 1678 1998 as a reimplementation of the program mfold with im-
 1679 proved thermodynamic parameters. In its initial version,
 1680 four major changes were made in mfold: (1) an improve-
 1681 ment on the methods for forcing base pairs; (2) a filter that
 1682 removed isolated WC or wobble base pairs has been added;
 1683 (3) the energy parameter for interior, internal and hairpin
 1684 loops were incorporated; (4) a new model for coaxial stack-
 1685 ing of helices. It predicts the lowest free energy structure
 1686 and a set of low energy structures. The new implemen-
 1687 tation also provided a user-friendly graphical interface for
 1688 Windows operating system. Subsequently, the first imple-
 1689 mentation was extended to include biomolecular folding; an
 1690 algorithm that finds low free energy structures common to
 1691 two sequences; the partition function algorithm and all free

1692 energy structures, and the constraints with enzymatic data
 1693 and chemical mapping data. The recent version includes
 1694 the partition function computation for secondary structures
 1695 common to two sequences and can perform stochastic sam-
 1696 pling of common structures [75]. Additionally, it contains
 1697 MaxExpect, which finds maximum expected accuracy struc-
 1698 tures [116], and a method for removal of pseudoknots, leav-
 1699 ing behind the lowest free energy pseudoknot-free struc-
 1700 ture.

- 1701 • **RNAfold** [80, 112]: It is one of the most used and efficient
 1702 folding tools. It computes the MFE secondary structure using
 1703 an efficient DP scheme and backtraces an optimum struc-
 1704 ture. It also allows computing the partition function using
 1705 McCaskill's algorithm, the matrix of base pairing probabili-
 1706 ties, and the centroid structure. It is part of the ViennaRNA
 1707 Package. Since its first version, it aims at suggesting an effi-
 1708 cient implementation of Zucker's algorithm with more flex-
 1709 ibility on the folding constraints. Many other versions have
 1710 been released, including a graphics processing unit (GPU)
 1711 implementation. The latest stable release of the ViennaRNA
 1712 Package is Version 2.5.0.
- 1713 • **LinearFold** [84]: For many decades, the DP techniques have
 1714 been the most accurate and fast at predicting pseudoknot-
 1715 free structure for short input RNA sequences. But for long
 1716 sequences, the prediction remains challenging because of
 1717 the computational time and the lack of accurate thermo-
 1718 dynamic energy parameters. In contrast to traditional DP
 1719 methods which are often bottom-up, LinearFold is a left-to-
 1720 right DP. The left-to-right DP consists of scanning the input
 1721 RNA sequence ϕ from left to right, maintaining a *stack* along
 1722 the way and performing one of the three actions (*push*,
 1723 *skip* or *pop*). The *stack* consists of a list of unpaired opening
 1724 bracket positions and at each position $j = 1 \dots L$, the three
 1725 actions consist respectively of 1) *push*: opening a bracket at
 1726 position j , 2) *skip*: unpaired nucleotide at position j and 3)
 1727 *pop*: closing the bracket at position j . Initially, LinearFold's
 1728 computational time was similar to the classical DP ($O(L^3)$)
 1729 because of the *pop* action that involves three free indices
 1730 (i.e. unpaired positions). But using a beam search heuris-
 1731 tic, the time complexity was then reduced to $O(Lb \log b)$,
 1732 where b is the beam size. The beam search is a popular
 1733 heuristic technique used in computational linguistics. This

1734 technique allows keeping only the top b highest-scoring (or
1735 low energy) states for each prefix of the input sequences.

1736 Although the score-based approaches for **RNA** structure pre-
1737 diction often offer good accuracy and generalization, the non-
1738 availability of the thermodynamic energy parameters for specific
1739 loops of extended sizes presents the main challenge for predict-
1740 ing long sequences (i.e. $L \geq 1,000$ nucleotides). Early **ML**-based
1741 methods aim to improve the energy parameters by learning the
1742 underlying folding patterns from a more considerable amount of
1743 training data. In the next section of this chapter, we will present
1744 some of the recent improvements in structure prediction using
1745 **ML**-based methods.

1746 2.1.2 *ML-based methods*

1747 In the previous section, we reviewed the score-based **RNA** sec-
1748 ondary structure prediction methods in general and four tools in
1749 particular, i.e. **Unfold**, **RNAstructure**, **RNAfold**, and **LinearFold**.
1750 These methods are thermodynamic methods that usually rely on
1751 experimentally energy parameters. For example, most experimen-
1752 tal energy parameters are available only for short **RNA** sequences
1753 (e.g. with a length of fewer than 200 nucleotides). This limitation
1754 significantly degrades the prediction performance of thermody-
1755 namic methods for long RNA sequences. In an attempt to improve
1756 these methods, **ML** methods have been proposed. This section
1757 presents an overview of existing **ML** methods, especially those
1758 used in [Chapter 3](#) for benchmark comparison with our proposed
1759 method.

1760 The **ML**-based methods for **RNA** secondary structure prediction
1761 can generally be classified into three categories according to **ML**'s
1762 subprocess, i.e., score scheme based on **ML**, preprocessing and
1763 postprocessing based on **ML**, and prediction process on **ML**. All
1764 the **ML**-based methods in these three categories trained their
1765 models in a supervised way [225].

1766 When using a scoring scheme based on **ML**, the parameter es-
1767 timation in the scoring scheme is first optimized using an **ML**
1768 model. The estimated parameters are then used to evaluate the
1769 scores of possible conformations. Difference scoring schemes can
1770 be refined by using that approach: the free energy parameters,
1771 weights, and probabilities. The free energy parameter-refining
1772 is the most popular because several thermodynamic parameters
1773 of the **NN** model have to be based on a large number of optimal

melting experiments and the experiments are time and labour-consuming. In fact, not all free energy changes in structural elements can be experimentally measured because of technical difficulties. Instead of refining the free energy parameters, some ML-based approaches scream through existing data of RNA structures to extract weights that consist of different features of RNA structure elements. These weights can be used as a scoring function for DP techniques. The advantage of such a scoring function is that it decouples structure prediction and energy estimation. However, learned weights have no explanations because of the ML black box.

Another alternative for predicting RNA structures is the stochastic context-free grammar (SCFG) [41, 102, 103, 151, 157, 214]. SCFGs allow building grammar rules and induce a joint probability distribution over possible RNA structure for a given sequence ϕ . In addition, the SCFG models specify probability parameters for each production rule in the grammar, which allow assigning a probability to each sequence generated by the grammar. These probability parameters are learned from datasets of RNA sequences associated with known secondary structures without carrying any external laboratory experiments [41].

Besides the ML-based methods that focus on refining the folding parameters, there are preprocessing and post-processing based on ML [77, 83, 226] and direct predicting process based on ML [111, 184, 187]. Preprocessing and postprocessing models allow for choosing the appropriate prediction method or set of prediction parameter sets and provide a means of determining the most likely structures among the possible outcomes that are useful for decision. The preprocessing and postprocessing ML tools are often based on a support vector machine (SVM).

Finally, it is possible to use ML techniques to predict RNA secondary structure directly or combine it with other algorithms in an end-to-end fashion. Below are some of the most used and recent ML-based tools for RNA secondary structure prediction.

- ContraFold[38]: Using the so-called probabilistic model, the conditional log-linear model (CLLM), ContraFold appeared for the first time in early 2006. It was the first probabilistic prediction tool outperforming the existing tools, including thermodynamic tools such as RNAfold and mfold. The CLLM is a flexible class of probabilistic models that generalizes upon SCFGs, using discriminative training and feature-rich scoring. The tool implements a CLLM incorporating most of the features found in typical thermodynamic

models allowing the tool to achieve the highest single sequence prediction accuracy to date when compared with the currently available probabilistic models.

- ContextFold [220]: In contrast to ContraFold, ContextFold utilizes a weighted approach based on ML. In particular, it uses a discriminative structured-prediction learning framework combined with an online learning algorithm. ContextFold uses a large training dataset of RNA sequences annotated with their corresponding structures to obtain an ML model made of 70,000 free parameters, which has several orders of magnitudes compared to traditional models (i.e. thermodynamic free energy parameters). At its first apparition, ContextFold’s model succeeded at the error reduction of about 50%. Still, some overfitting concerns have been reported when using the tool, especially for predicting structures with large unpaired regions.
- Mfold2 [160]: It is one of the most recent ML-based tools for predicting the secondary structure of RNA molecules. Its particularity is the ML technique used, a ML it also belongs to the weighted approach based on ML since the resulting model of a deep neural network (DNN) is a set of weight parameters. Mfold2’s DNN uses the max-margin framework with thermodynamic regularization. It made the folding scores predicted by Mfold2 and the free energy calculated by the thermodynamic parameters as close as possible. This method has shown robust prediction on both sequences and families of natural RNAs, suggesting that the weighted ML approaches can compensate for the gaps in the thermodynamic parameter approaches.

Although ML methods provide substantial improvements compared to traditional methods such as thermodynamic and comparative sequence analysis [161, 175], they often lack physical principles (training data are mostly obtained through phylogenetic analyses) and present some over-fitting concerns [152]. In addition to the over-fitting problems partially due to few data availability, ML methods do not provide dynamic information on RNA folding for the same reason. In Chapter 3, we will introduce our approach that aims at predicting an ensemble structure, which allows us to derive some dynamic information and contrasts the methods previously presented.

1857 2.1.3 *Prediction tools for pseudoknotted RNA sequences*

1858 In the introduction, we have provided the importance of pseudo-
 1859 knot interaction in realizing biological functions, and different
 1860 pseudoknot patterns have been reviewed. This section introduces
 1861 a couple of tools for predicting RNA pseudoknotted structures
 1862 that will be used in the benchmark results presented in Chapter 5.

1863 Folding RNA sequences with pseudoknotted interactions is
 1864 computationally more expensive than a pseudoknot-free tar-
 1865 get. Specifically, the time complexity of the pseudoknot-free sec-
 1866 ondary structure prediction is $O(L^3)$ when using dynamic pro-
 1867 gramming approaches such as RNAfold, or less with heuristic
 1868 folding methods (e.g. $O(L)$ for LinearFold and $O(L^2 \log L)$). By
 1869 contrast, when considering a special class of pseudoknots, the
 1870 time complexity of folding goes up to $O(L^6)$ for an exact ther-
 1871 modynamic prediction using a dynamic programming approach
 1872 such as [150]. When Using heuristic methods, the time complex-
 1873 ity slows down to $O(L^4)$ (e.g. tools such as IPknot and HotKnots)
 1874 or $O(L^3)$ for tool such as HFold.

- 1875 • pKiss [90]: The program pKiss appears the first time in
 1876 2014 as an updated version of the program pknotsRG[144]
 1877 which is a module of the RNA abstract shapes analysis
 1878 RNAshapes [90]. Initially, the program pknotsRG was built
 1879 for the prediction of some special class of pseudoknots (un-
 1880 knotted structures and H-type pseudoknots). Later on, it
 1881 was extended to predict RNA structures that exhibit kissing
 1882 hairpin motifs in an arbitrarily nested fashion, requiring
 1883 $O(L^4)$ time. In addition to predicting the kissing hairpin mot-
 1884 tifs, pKiss also provides new features such as shape analy-
 1885 sis, computation of probabilities, different folding strategies
 1886 and different dangling base models.
- 1887 • IPknot [162]: it was first introduced in a paper by Kengo
 1888 and his collaborators in 2011 as a novel computational tool
 1889 for predicting RNA secondary structure with pseudoknots
 1890 using integer programming technique. IPknot uses the
 1891 maximum expected accuracy (MEA) as a scoring function,
 1892 and the maximizing expected accuracy problem is solved
 1893 using integer programming with threshold cut. IPknot de-
 1894 composes a pseudoknotted structure into a set of pseudoknot-
 1895 free substructures and approximates a base-pairing prob-
 1896 ability distribution that considers pseudoknots, leading
 1897 to the capability of modelling a comprehensive class of

1898 pseudoknots and running quite fast. In addition to single
1899 sequence analysis, IPknot can also predict the consensus
1900 secondary structure with pseudoknots when a multiple
1901 sequence alignment is given.

- 1902 • HotKnots [147]. In contrast to the previously mentioned
1903 tools, HotKnots implements a heuristic algorithm based on
1904 the simple idea of iteratively forming stable stems. The algo-
1905 rithm explores many alternative secondary structures using
1906 a free energy minimization for pseudoknot-free secondary
1907 structures. Several other additions of a single substructure
1908 are considered for each structure formed at each step, re-
1909 sulting in a tree of candidate structures. The criterion for
1910 determining which substructures to add to partially formed
1911 structures at successive levels of the tree was also new. Sim-
1912 ilar to previous algorithms, energetically favourable sub-
1913 structures called *hotspots* are found by a call to Zuker's
1914 algorithm, with the constraint that no base already paired
1915 may be in the structure.

1916 Despite the higher computational complexity of pseudoknots,
1917 it is still important to account for them as they occur in natu-
1918 ral RNA and are relevant for RNA function. We have reviewed
1919 three mainly used RNA secondary prediction tools (pKiss, IPknot,
1920 HotKnots) that support the two pseudoknot patterns (i.e. the H-
1921 type and K-type) considered in Chapter 5. In addition to the
1922 computational complexity, existing methods lack experimentally
1923 measured energy parameters for pseudoknot interactions. There-
1924 fore, they mostly rely or do not on approximated energy param-
1925 eters, which may influence the predictions. Only IPknot and
1926 HotKnots will be used among these tools when designing pseu-
1927 doknotted RNA structures. HotKnots predicts the free energy of
1928 pseudoknotted structure based on recently updated energy pa-
1929 rameters, whereas IPknot does not.

1930 So far, we have presented tools that predict a single stable and
1931 static RNA secondary structure for a given RNA sequence, includ-
1932 ing pseudoknots or not. More often than not, the ncRNA functions
1933 are associated with the RNAs' ability to undergo specific confor-
1934 mational changes, as is the case for riboswitches. The function of
1935 an RNA molecule thus is usually poorly described by its ground
1936 state structure and instead has to be studied as a dynamic en-
1937 semble of structures [35, 137]. The following section will review
1938 some computational methods that address the folding dynamics
1939 of RNA molecules.

1940 2.2 RNA KINETICS

1941 The previous section introduced how pseudoknot-free secondary
 1942 structures with their thermodynamic properties can be predicted.
 1943 It also introduced some statistical methods that do not only rely
 1944 on the thermodynamic principle but training data obtained from
 1945 phylogenetic analysis, mainly the ML methods. However, the
 1946 methods used for predictions do not tell us anything about how
 1947 the structures change over time and how they are related to
 1948 each other. This section discusses the folding dynamics of RNA
 1949 molecules.

1950 The folding of RNA molecules is remarkably more complex. It
 1951 is a result of the delicate balance between multiple factors: the
 1952 chain entropy, ion-mediated electrostatic interactions and solva-
 1953 tion effect, base pairing and stacking, and other non-canonical
 1954 interactions [22]. It is a dynamic process governed by a constant
 1955 formation or dissolving of base pairs. In other terms, the RNA
 1956 molecule navigates its structure space by following a free energy
 1957 landscape. Here, the free energy landscape is a high-dimensional
 1958 space of all possible secondary structures (Σ_ϕ) weighted by their
 1959 free energy ΔG .

1960 As usually done, the kinetics is modelled as a continuous-time
 1961 Markov chain [114], where populations of structure evolve accord-
 1962 ing to transition rates. In this context, an Arrhenius formulation
 1963 is commonly used to derive elementary transition from state i to
 1964 state j ; where $\Delta G_{i \rightarrow j}^\ddagger$ is the activation barrier separating i from j ,
 1965 and $\beta = 1/k_B T$ is the inverse thermal energy (mol/kcal).

$$k_{i \rightarrow j} = k_0 \exp(-\beta \Delta G_{i \rightarrow j}^\ddagger) \quad (2.2)$$

1966 Here k_0 is the actual rate constant, solvent-dependent. Three rate
 1967 models describing elementary steps in the structure space are
 1968 often used to study RNA folding dynamics:

- 1969 1. The base stack model [221–223]: it uses base stacks as ele-
 1970 mentary kinetic move. A move consists of an addition or a
 1971 breaking of a base stack with $\Delta G_{i \rightarrow j}^\ddagger$ equal to the change in
 1972 the entropic free energy $T\Delta S$ and the enthalpy ΔH , respec-
 1973 tively.
- 1974 2. The base pair model [24, 53]: it uses base pair as elementary
 1975 kinetic steps which gives the finest resolution, but at the
 1976 cost of computation time. Here $\Delta G_{i \rightarrow j}^\ddagger = \Delta G/2$ where ΔG is

1977 the energy change from state i to state j or $\Delta G_{i \rightarrow j}^\ddagger = \Delta G$ for
 1978 $\Delta G \geq 0$.

1979 3. The helix stem model [87, 121]: the elementary move is the
 1980 creation or deletion of a helix stem. It provides a coarse-
 1981 grained description of the dynamics where free energy
 1982 changes ($\Delta G_{i \rightarrow j}^\ddagger$) due to stem formation guiding the folding
 1983 process.

1984 The different rate models can lead to different folding path-
 1985 ways. The key factor that distinguishes the different rate mod-
 1986 els is whether the barrier is determined by $(\Delta H, \Delta S)$ or by ΔG .
 1987 The $(\Delta H, \Delta S)$ values for different RNA base stacks show well-
 1988 separated discrete hierarchies, whereas the ΔG values show no
 1989 such large separation. For two typical base stacks, 5'AU-AU3' and
 1990 5'UC-GA3', the difference $\Delta(\Delta H_{stack}, \Delta S_{stack}) = (7.4 \text{ kcal/mol},$
 1991 $20 \text{ kcal/mol})$ is much larger than the difference $\Delta(\Delta G_{stack}) =$
 1992 1.4 kcal/mol [168]. Because of this fact, different models can give
 1993 different folding kinetics.

1994 Depending on the rate model used, the following master-equation
 1995 describe the population kinetics $p_i(t)$ for the i^{th} state ($i = 1 \dots \Omega$,
 1996 where Ω is the total number of chain conformations).

$$\frac{dp_i(t)}{dt} = \sum_{j \in \Omega} k_{j \rightarrow i} p_j(t) - k_{i \rightarrow j} p_i(t), \quad (2.3)$$

where $k_{j \rightarrow i}$ and $k_{i \rightarrow j}$ are the rate constants for the respective transitions. The equivalent matrix form of Equation 2.3 is given by:

$$\frac{d\mathbf{p}(t)}{dt} = \mathbf{M} \cdot \mathbf{p}, \quad (2.4)$$

1997 where $\mathbf{p} = (p_1, \dots, p_\Omega)$ is a column vector representing the fre-
 1998 quency of structure at state (i, \dots, Ω) and, \mathbf{M} is the rate matrix
 1999 defined as:

$$\mathbf{M}_{ij} = \begin{cases} k_{i \rightarrow j}, & \text{if } i \neq j \\ -\sum_{j \neq i} k_{ij}, & \text{if } i = j \end{cases} \quad (2.5)$$

2000 For a given initial folding condition $p_i(0)$, the Equation 2.4 is
 2001 solvable by diagonalizing the rate matrix \mathbf{M} and, the solution is
 2002 the population kinetics $\mathbf{p}(t)$ for $t > 0$ is given by:

$$\mathbf{p}(t) = \sum_{m=1}^{\Omega} C_m \mathbf{n}_m \exp -\lambda_m t \quad (2.6)$$

2003 where $-\lambda$ and \mathbf{n}_m are the m^{th} eigenvalue and eigenvector of
 2004 the rate matrix \mathbf{M} , and C_m is the coefficient that is dependent on
 2005 the initial condition. The eigenvalue spectrum gives the rates of
 2006 the kinetic modes of the system.

2007 Simulating the RNA dynamics using [Equation 2.3](#) has some
 2008 limitations. The solution to the master-equation given by [Equa-](#)
 2009 [tion 2.6](#) can only give ensemble-average macroscopic kinetics
 2010 and cannot give detailed information about the microscopic path-
 2011 ways [224]. Moreover, the number of structures (Ω) increases
 2012 rapidly with the RNA sequence length L . Therefore, the master
 2013 equation is often limited to short RNA sequences. Because of these
 2014 limitations, kinetics-cluster methods are alternatively used. The
 2015 basic idea of the kinetic-cluster method is to classify the large
 2016 structural ensemble into a much-reduced system of clusters (of
 2017 macrostates) such that the inter-cluster transitions can represent
 2018 the overall kinetics. Although both the master-equation and the
 2019 kinetic-cluster methods can predict the macroscopic kinetics, the
 2020 kinetic-cluster approach has the unique advantage of providing
 2021 direct information on the microscopic pathway statistics from
 2022 the inter-cluster transitions [224]. Both approaches are based on
 2023 the complete conformational ensemble. An alternative approach,
 2024 implemented in [kinwalker](#) [62], used the observation that folded
 2025 intermediates are generally locally optimal conformations. Like
 2026 thermodynamic methods for static RNA secondary structure pre-
 2027 diction, experimental studies usually play an essential role in
 2028 guiding computational methods in studying RNA folding dynam-
 2029 ics. Several recent observations are discussed in the following
 2030 paragraph.

2031 In folding experiments, Pan and coworkers observed two kinds
 2032 of pathways in the free energy landscape of a natural ribozyme
 2033 [139]. Firstly, the investigations revealed fast-folding pathways, in
 2034 which a subpopulation of RNAs folded rapidly into the native state.
 2035 However, the second population quickly reached metastable mis-
 2036 folded states, then slowly folded into the native structure. In some
 2037 cases, these metastable states are functional. These phenomena
 2038 are direct consequences of the rugged nature of the RNA folding
 2039 landscape [180].

2040 The experiments performed by Russell and coworkers also
 2041 revealed the presence of multiple deep channels separated by
 2042 high energy barriers on the folding landscape, leading to fast
 2043 and slow folding pathways [156]. The formal description of the
 2044 above mechanism, called the kinetic partitioning mechanism,
 2045 was first introduced by Guo and Thirumalai in the context of

2046 protein folding [69]. These metastable conformations constitute
2047 competing attraction basins in the free energy landscape where
2048 RNA molecules are temporarily trapped. However, *in vivo*, folding
2049 into the native states can be promoted by molecular chaperones
2050 [20], which means that the active structure depends on factors
2051 other than the sequence. This may raise some discrepancies when
2052 comparing thermodynamic modelling to actual data.

2053 The experimental verification of the rate model is also a chal-
2054 lenge because the microscopic elementary processes are hidden in
2055 the ensemble averages of the measured kinetics. Many researchers
2056 believe that single-molecule experiments may provide a discern-
2057 ing measure with careful extrapolation to the force-free case.
2058 All atom-simulations with a reliable force field and sampling
2059 method are highly valuable for providing detailed atomistic con-
2060 figurations for the transition state [22]. Alternatively, systematic
2061 theory-experiment tests as done in [224] for designed sequences
2062 can also provide critical assessment for the different rate models.

2063 In sum, studying the folding of RNA molecules as a dynamic
2064 ensemble of structures is of central importance in describing their
2065 functions, and experimental observations often guide the compu-
2066 tational methods. Some of the recent experimental observations
2067 have been reviewed in this section. Among them, the kinetic par-
2068 titioning mechanism is of interest in this work. It revealed the
2069 presence of multiple deep channels separated by high-energy
2070 barriers on the folding landscape, which leads to fast and slow
2071 folding pathways. The folding tool we suggest in Chapter 3 is in-
2072 spired by this mechanism and predicts fat RNA folding pathways.
2073 The predicted pathways, therefore, allow us to derive dynamic
2074 information on RNA folding.

2075 2.3 CONCLUSION

2076 In this chapter, we have presented the RNA folding in two main
2077 steps: (1) the prediction of the secondary structure of RNA, which
2078 represents the static part of the folding process; (2) the RNA ki-
2079 netics, which aim at modelling the dynamics of the folding. RNA
2080 secondary structure prediction was introduced as an optimiza-
2081 tion problem, and a review of existing methods and tools was
2082 presented. Of particular importance in this thesis's context is
2083 that existing tools for predicting RNA secondary structures of-
2084 ten present some limits in computational time for longer RNA
2085 sequences. Mainly the existing tools do not give dynamical infor-
2086 mation, as few data are available on structural dynamics. Sim-

2087 ulating the folding kinetics of long RNA molecules is also of an
2088 essential limit because it requires a full enumeration of the struc-
2089 tural space in most cases. In the next chapter, we will present our
2090 thesis's first result, which aims to predict RNA folding pathways
2091 efficiently using the FFT. The predicted pathways allow us to de-
2092 rive energetically suboptimal structures from which we model
2093 the RNA folding kinetics with fewer secondary structures.

2094

2095 RAFFT: EFFICIENT PREDICTION OF
2096 FAST-FOLDING PATHWAYS OF RNAs

2097 This chapter introduces a novel heuristic algorithm to predict
2098 an ensemble of metastable RNA secondary structures for a given
2099 sequence ϕ . The algorithm is inspired by the kinetic partitioning
2100 mechanism, by which molecules follow alternative folding
2101 pathways to their native structure, some much faster than others.
2102 Similarly, our algorithm RAFFT generates an ensemble of concurrent
2103 folding pathways ending in multiple metastable structures
2104 for each given sequence. We then use the ensemble structures
2105 as finite ensemble states in which the RNA sequence can be at a
2106 given time, and the energy difference from one state to another is
2107 then used to derive a stem rate model. Therefore, our algorithm
2108 also acts as a folding kinetic ansatz. Much of the material in this
2109 chapter has been previously described in [138].

2110 3.1 MATERIAL AND METHODS

2111 The computational time is one of the challenges for the existing
2112 tool in folding long RNA molecules. The method we present
2113 in this work aims to improve the existing RNA folding tools re-
2114 viewed in Chapter 2. It is based on the FFT and inspired by the
2115 kinetic partitioning mechanism. As presented in Chapter 1, the
2116 FFT allows reducing the computational time of the correlation
2117 between two sequences. We use the same ideal in the context of
2118 this work to faster predict RNA folding pathways by analyzing
2119 high correlation positional lag between an RNA sequence and its
2120 complementary copy, especially for longer sequences. We, there-
2121 fore, derive a kinetics ansatz from the structural ensemble of the
2122 predicted folding paths. This section describes our RNA pathways
2123 prediction method and the kinetics ansatz derived from the pre-
2124 dicted structural ensemble. In addition, it provides a description
2125 of the benchmark dataset used to assess our method performance
2126 and the comparison protocols.

²¹²⁷ 3.1.1 *RAFFT's algorithm description*

RAFFT starts from a sequence of nucleotides $\phi = (\phi_1 \dots \phi_L)$ of length L , and its associated unfolded structure σ . We first create a numerical representation of ϕ where each nucleotide is replaced by a unit vector of 4 components:

$$A \rightarrow \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}, U \rightarrow \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix}, C \rightarrow \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix}, G \rightarrow \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} \quad (3.1)$$

This encoding gives us a $(4 \times L)$ -matrix we call X , where each row corresponds to a nucleotide as shown below:

$$X = \begin{pmatrix} X^A \\ X^C \\ X^G \\ X^U \end{pmatrix} = \begin{pmatrix} X^A(1) & X^A(2) & \dots & X^A(L) \\ X^C(1) & X^C(2) & \dots & X^C(L) \\ X^G(1) & X^G(2) & \dots & X^G(L) \\ X^U(1) & X^U(2) & \dots & X^U(L) \end{pmatrix} \quad (3.2)$$

For example, $X^A(i) = 1$ if $\phi_i = A$. Next, we create a second copy $\bar{\phi} = (\bar{\phi}_L \dots \bar{\phi}_1)$ for which we reversed the sequence order. Then, each nucleotide of $\bar{\phi}$ is replaced by one of the following unit vectors:

$$\bar{A} \rightarrow \begin{pmatrix} 0 \\ 0 \\ 0 \\ w_{AU} \end{pmatrix}, \bar{U} \rightarrow \begin{pmatrix} w_{GU} \\ 0 \\ 0 \\ 0 \end{pmatrix}, \bar{C} \rightarrow \begin{pmatrix} 0 \\ 0 \\ w_{GC} \\ 0 \end{pmatrix}, \bar{G} \rightarrow \begin{pmatrix} 0 \\ w_{GC} \\ 0 \\ w_{GU} \end{pmatrix}. \quad (3.3)$$

²¹²⁸ \bar{A} (respectively $\bar{U}, \bar{C}, \bar{G}$) is the complementary of A (respectively ²¹²⁹ U, C, G). w_{AU}, w_{GC}, w_{GU} represent the weights associated with ²¹³⁰ each canonical base pair, and they are chosen empirically. We call ²¹³¹ this complementary copy \bar{X} , the mirror of X .

To search for stems, we use the complementary relation between X and \bar{X} with the correlation function $\text{cor}(k)$. This correlation is defined as the sum of individual X and \bar{X} row correlations:

$$\text{cor}(k) = \sum_{\alpha \in \{A, U, C, G\}} c_{X^\alpha, \bar{X}^\alpha}(k), \quad (3.4)$$

where a row correlation between X and \bar{X} is given by:

$$c_{X^\alpha, \bar{X}^\alpha}(k) = \sum_{\substack{1 \leq i \leq L \\ 1 \leq i+k \leq L}} \frac{X^\alpha(i)\bar{X}^\alpha(i+k)}{\min(k, 2L-k)}. \quad (3.5)$$

2132 For each $\alpha \in \{A, U, C, G\}$, $X^\alpha(i) \times \bar{X}^\alpha(i+k)$ is non zero if sites
 2133 i and $i+k$ can form a base pair, and will have the value of the
 2134 chosen weight as described above. If all the weights are set to
 2135 1, $\text{cor}(k)$ gives the frequency of base pairs for a positional lag k .
 2136 Although the correlation naively requires $O(L^2)$ operations, it
 2137 can take advantage of the [FFT](#) which reduces its complexity to
 2138 $O(L \log(L))$.

2139 Large $\text{cor}(k)$ values between the two copies indicate positional
 2140 lags k where the frequency of base pairs is likely to be high. How-
 2141 ever, this does not allow to determine the exact stem positions.
 2142 Hence, we use a sliding window strategy to search for the largest
 2143 stem within the positional lag (since the copies are symmetrical,
 2144 we only need to slide over one-half of the positional lag). Once
 2145 the largest stem is identified, we compute the free energy change
 2146 associated with the formation of that stem. Next, we perform the
 2147 same search for the n highest correlation values, which gives us
 2148 n potential stems. Then, we define as the current structure the
 2149 stem with the lowest free energy. Here, free energies were com-
 2150 puted using Turner2004 energy parameters through ViennaRNA
 2151 package [API](#) [112].

2152 We are now left with two independent parts, the interior and
 2153 the exterior of the newly formed stem. If the exterior part is com-
 2154 posed of two fragments, they are concatenated into one. Then,
 2155 we apply recursively the same procedure on the two parts in-
 2156 dependently in a *breadth-first* fashion to form new consecutive
 2157 base pairs. The procedure stops when no base pair formation
 2158 can improve the energy. When multiple stems can be formed in
 2159 these independent fragments, we combine all of them and pick
 2160 the composition with the best overall stability. If too many com-
 2161 positions can be formed, we restrict this to the 10^4 bests in terms
 2162 of energy. [Figure 3.1](#) shows an example of execution to illustrate
 2163 the procedure.

2164 The algorithm described so far tends to be stuck in the first
 2165 local minima found along the folding trajectory. To alleviate this,
 2166 we implemented a stacking procedure where the N best trajec-
 2167 tories are stored in a stack and evolved in parallel. Like the initial
 2168 version, the algorithm starts with the unfolded structure; then,
 2169 the N best potential stems are stored in the first stack. From these
 2170 N structures, the procedure tries to add stems in the unpaired
 2171 regions left and saves the N best structures formed. Once no stem
 2172 can be formed, the algorithm stops and output the structure with
 2173 the best energy found among the structures stored in the last
 2174 stack. This algorithm leads to the construction of a graph we call

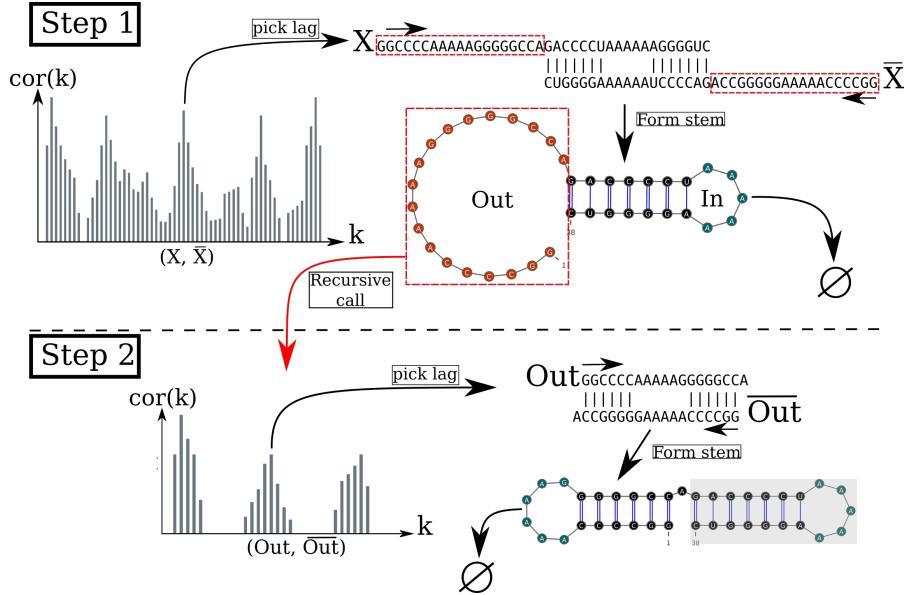


Figure 3.1: Algorithm execution for one example sequence which requires two steps. (Step 1) From the correlation $cor(k)$, we select one peak which corresponds to a position lag k . Then, we search for the largest stem and form it. Two fragments, “In” (the interior part of the stem) and “Out” (the exterior part of the stem), are left, but only the “Out” may contain a new stem to add. (Step 2) The procedure is called recursively on the “Out” sequence fragment only. The correlation $cor(k)$ between the “Out” fragment and its mirror is then computed and analyzing the k positional lags allows to form a new stem. Finally, no more stem can be formed on the fragment left (colored in blue), so the procedure stops.

2175 a *fast-folding graph*. In this graph, two structures are connected if
 2176 the transition from one to another corresponds to the formation
 2177 of a stem or if the two structures are identical. [Figure 3.2](#) shows
 2178 an example of a *fast-folding graph* produced by RAFFT for $N = 5$.

2179 This section presented the complete procedure implemented in
 2180 our proposed tool RAFFT. The procedure resulted in an ensemble
 2181 of concurrent folding pathways ending in multiple metastable
 2182 secondary structures. The connections in each folding pathway
 2183 are dictated by the formation of stems, resulting in an energy
 2184 increase. The different folding pathways connected to the initial
 2185 unfolded structure form a fast folding graph. The ensemble of
 2186 secondary structures constituting the fast folding graph is then
 2187 used to build our kinetics ansatz where the transitions follow
 2188 the Metropolis rules, i.e. no barriers between structures. The

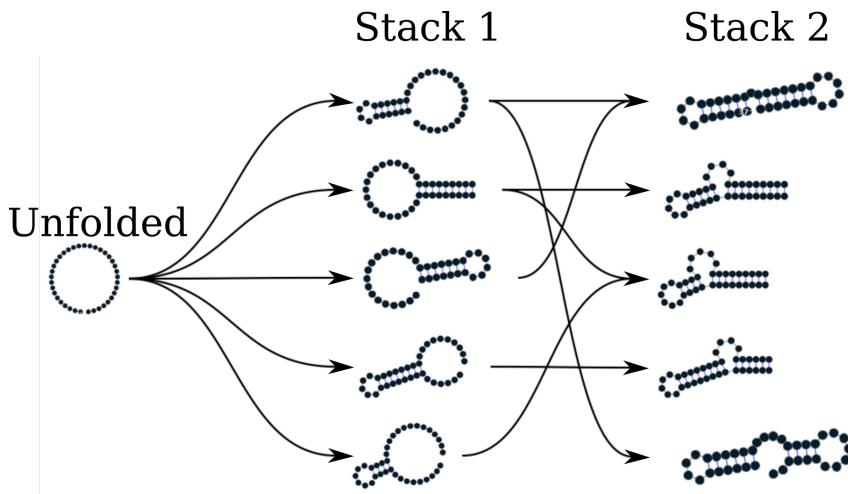


Figure 3.2: Fast folding graph constructed using RAFFT. In this example, the sequence is folded in two steps. The algorithm starts with the unfolded structure on the left. The $N = 5$ best stems are stored in stack 1. From stack 1, multiple stems formation are considered, but only the $N = 5$ best are stored in stack 2. Structures are ordered (from top to bottom) by energy in each stack. All secondary structure visualizations were obtained using VARNA [31].

2189 following section provides more details on our proposed kinetics
2190 ansatz.

2191 3.1.2 Kinetic ansatz

2192 Now that the RNA pathway prediction algorithm is described,
2193 we provide the ingredients needed to extract dynamic folding
2194 information from the previously generated fast folding graph in
2195 this section.

2196 The folding kinetic ansatz used here is derived from the fast-
2197 folding graph and allows us to model the slow processes in RNA
2198 folding. As described in Figure 3.2, transitions can occur from
2199 left to right (and right to left) but not vertically. The fast-folding
2200 graph follows the idea that parallel pathways quickly reach their
2201 endpoints; however, when the endpoints are non-native states,
2202 this ansatz allows slowly folding back into the native state [139].

2203 Using the master-equation (See Equation 2.3), the traditional
2204 kinetic approach often starts by enumerating the whole space (or
2205 a carefully chosen subspace) of structures using RNAsubopt. Next,
2206 this ensemble is divided into local attraction basins separated
2207 from one another by energy barriers. This coarsening is usually

2208 done with the tool called **barriers**. Then, following the Arrhenius
 2209 formulation (See [Equation 2.2](#)) , one simulates a coarse grained
 2210 kinetics between basins.

In contrast to traditional kinetics approaches, the connected structures in the RAFFT's fast-folding graph are not always separated by activation barrier energies. Therefore, we computed the transition rates $k_{i \rightarrow j}$ using the Metropolis [[101](#)] formulation defined as follow

$$k_{i \rightarrow j} = \begin{cases} k_0 \times \min(1, \exp(-\beta \Delta(\Delta G_{i \rightarrow j}))), & \text{if } \sigma_i \in \mathcal{M}(\sigma_j) \\ 0, & \text{else} \end{cases}, \quad (3.6)$$

2211 where $\Delta \Delta G_{i \rightarrow j} = \Delta G_j - \Delta G_i$ is the free energy change between
 2212 structure σ_i and σ_j . Here, k_0 is a conversion constant that we
 2213 set to 1 for the sake of simplicity and we initialize the popula-
 2214 tion $p_i(0)$ with only unfolded structures; therefore, the trajectory
 2215 represents a complete folding process. The frequency of a struc-
 2216 ture σ_i evolves according to the master [Equation 2.3](#). Due to this
 2217 approximation, we referred to our approach as a *kinetic ansatz*

2218 In sum, based on the [FFT](#), we constructed a method that al-
 2219 lows generating an ensemble of secondary structures by a suc-
 2220 cessive formation of stems. Using this ensemble, we derived a
 2221 kinetics ansatz in which transitions between structures follow the
 2222 Metropolis rules. We assess the performance of our tool by com-
 2223 paring its predictions to existing tools using benchmark datasets.
 2224 The following section briefly describes the datasets used in this
 2225 work, including the clean procedure applied to the initial datasets.

2226 3.1.3 Benchmark datasets.

2227 Measuring the performance of computational RNA folding tools
 2228 can be quite a challenging task. A perfect validation procedure
 2229 will require a comparison to experimental data, which in practice
 2230 are not also perfect and are very expensive. In the context of this
 2231 work, we perform *in silico* validation using benchmark datasets,
 2232 which is a collection of native sequence structures. Because our
 2233 proposed method produces kinetics and static structure predic-
 2234 tions., we assess the performance of both tasks separately and
 2235 using a different dataset. This section presents the two datasets.

2236 To build the dataset for the folding task application, we started
 2237 from the ArchiveII dataset derived from multiple sources [[4, 9,](#)
 2238 [17, 29, 32, 61, 70, 71, 124, 158, 164, 176, 181, 183, 207, 229, 230](#)].
 2239 We first removed all the structures with pseudoknots, since the

2240 tools considered here do not handle these loops. Next, using the
 2241 Turner2004 energy parameters, we evaluated the structures' ener-
 2242 gies and removed all the unstable structures: structures with en-
 2243 ergies $\Delta G > 0$. This dataset is composed of 2,698 sequences with
 2244 their corresponding known structures. 240 sequences were found
 2245 multiple times (from 2 to 8 times); 19 of them were mapped to
 2246 different structures. For the sequences that appeared with differ-
 2247 ent structures, we picked the structure with the lowest energy. In
 2248 the end we arrived at a dataset with 2,296 sequences-structures.

2249 For the kinetics task, there is no existing standard procedure or
 2250 dataset allowing to validate or not a computational tool. However,
 2251 for the validation of our kinetic ansatz, we used the [CFSE RNA](#) se-
 2252 quence and classic bi-stable sequence **GGCCCCUUUUGGGGGCCA-**
 2253 **GACCCCUAAAGGGGUC**.

2254 In sum, two different dataset sets are used to assess RAFFT per-
 2255 formance: the first one, Archive II consists of 2,296 sequences-
 2256 structures used for the prediction task, and one which contains
 2257 two sequences, the CFSE and a bistable sequence for the kinetic
 2258 study. The following section describes the benchmarking proto-
 2259 cols for both tasks.

2260 3.1.4 Structure prediction protocols

2261 The static RNA structure prediction and the RNA kinetic perfor-
 2262 mances of our proposed tool RAFFT are evaluated separately. This
 2263 section describes the evaluation protocols for both performances
 2264 and the different tool parameters used throughout.

2265 To evaluate the structure prediction accuracy of the proposed
 2266 method, we compared RAFFT to five recent secondary structure
 2267 pseudoknot-free prediction tools. The five tools include [ML](#)-based
 2268 methods ([Mfold2 0.1.1](#) and [Contrafold](#)) and score-based meth-
 2269 ods ([RNAfold 2.4.13](#), [Linearfold](#), and [RNAsstructure](#)). To com-
 2270 pute the [MFE](#) structure for the score-based methods, we used the
 2271 default parameters and the Turner2004 set of energy parameters.
 2272 We also computed the [ML](#) predictions using the default param-
 2273 eters. Therefore, only one structure prediction per sequence for
 2274 these tools was used for the statistics.

2275 Two parameters are critical for RAFFT, the number of posi-
 2276 tional lags in which stems are searched, and the number of struc-
 2277 tures stored in the stack. For our computational experiments,
 2278 we searched for stems in the $n = 100$ best positional lags and
 2279 stored $N = 50$ structures. The correlation function $\text{cor}(k)$ which

2280 allows to choose the positional lags is computed using the weights
 2281 $w_{GC} = 3$, $w_{AU} = 2$, and $w_{GU} = 1$.

2282 To assess the performance of RAFFT, we analyzed the output in
 2283 two different ways. First, we considered only the structure with
 2284 the lowest energy found for each sequence. This procedure allows
 2285 us to assess RAFFT performance in predicting the MFE structure.
 2286 Second, we computed the accuracy of all $N = 50$ structures saved
 2287 in the last stack for each sequence and displayed only the best
 2288 structure in terms of accuracy. As mentioned previously in Chapter
 2289 2, the lowest energy structure found may not be the active
 2290 structure. Therefore, this second assessment procedure allows us
 2291 to show whether one of the pathways is biologically relevant.

2292 We used two metrics to measure the prediction accuracy: the
 2293 PPV and the sensitivity. The PPV measures the fraction of correct
 2294 base pairs in the predicted structure, while the sensitivity mea-
 2295 sure the fraction of base pairs in the accepted structure that are
 2296 predicted. These metrics are defined in Chapter 1 (See definitions
 2297 11, 12). To be consistent with previous studies, we computed these
 2298 metrics using the scorer tool provided by Matthews *et al.* [122],
 2299 which also provides a more flexible estimate where shifts are
 2300 allowed.

2301 Further more, we used a PCA to visualize the loop diversity in
 2302 the predicted structures for each folding tool considered here. To
 2303 extract the weights associated with each structure loop from the
 2304 dataset, we first converted the structures into weighted coarse-
 2305 grained tree representation [169]. In the tree representation, the
 2306 nodes are generally labelled as E (exterior loop), I (interior loop),
 2307 H (hairpin), B (bulge), S (stacks or stem-loop), M (multi-loop)
 2308 and R (root node). We separately extracted the corresponding
 2309 weights for each node, and the weights are summed up and
 2310 then normalized. Excluding the root node, we obtained a table
 2311 of 6 features and n entries. This allows us to compute a 6×6
 2312 correlation matrix that we diagonalize using the eigen routine
 2313 implemented in the scipy package. For visual convenience, the
 2314 structure compositions were projected onto the first two principal
 2315 components (PCs).

2316 Finally, the CFSE and a bistable RNA sequence are used to assess
 2317 the kinetic performance. For each sequence, initial conditions
 2318 are chosen for both Treekin and RAFFT to simulate the kinetic
 2319 trajectories. Both kinetics are simulated using the master equation
 2320 described in Chapter 2 (Equation 2.3) but with different transition
 2321 rules, Treekin uses the Arrhenius rules whereas RAFFT uses the

2322 Metropolis rules. The following section describes the statistical
2323 results obtained for both kinetics and structure prediction tasks.

2324 **3.2 EXPERIMENTAL RESULTS**

2325 The validation of our results is purely statistical, i.e. using statis-
2326 tical methods such as *t*-test and regression to compare different
2327 tool performance data. Based on the previously mentioned limi-
2328 tations of existing tools, we evaluate three main RAFFT potential
2329 improvements: the running time for the folding or pathways pre-
2330 diction, the quality of the predicted pathways and the RNA fold-
2331 ing kinetics. This section discusses each of those performances in
2332 comparison to the existing tools.

2333 **3.2.1 RAFFT's run time and scalability**

2334 The first input of our method is a potential improvement to the
2335 CPU time of existing tools. This section focuses on analyzing
2336 RAFFT's running time compared to existing methods. Four dif-
2337 ferent tools are considered: RNAfold, ContraFold, RNAstructure
2338 and LinearFold. All of them are MFE estimates implementing a DP
2339 with cubic time complexity($O(L^3)$), except for ContraFold which
2340 implements a ML approach. When using the heuristic implemen-
2341 tation of LinearFold, the time complexity is linear while losing
2342 the MFE estimation. We will first discuss RAFFT theoretical time
2343 complexity before comparing the empirical execution times to
2344 the existing tools.

2345 The complexity of RAFFT's algorithm depends on the number
2346 and size of the stems formed. The main operations performed
2347 for each stem formed are: (1) the evaluation of the correlation
2348 function $\text{cor}(k)$, (2) the sliding-window search for stems, and (3)
2349 the energy evaluation. We based our approximate complexity on
2350 the correlation evaluation since it is the more computationally
2351 demanding step; the other operations only contribute a multi-
2352 plicative constant at most. The best case is the trivial structure
2353 composed of one large stem where the algorithm stops after eval-
2354 uating the correlation on the complete sequence. At the other
2355 extreme, the worst case is one where at most $L/2$ stems of size 1

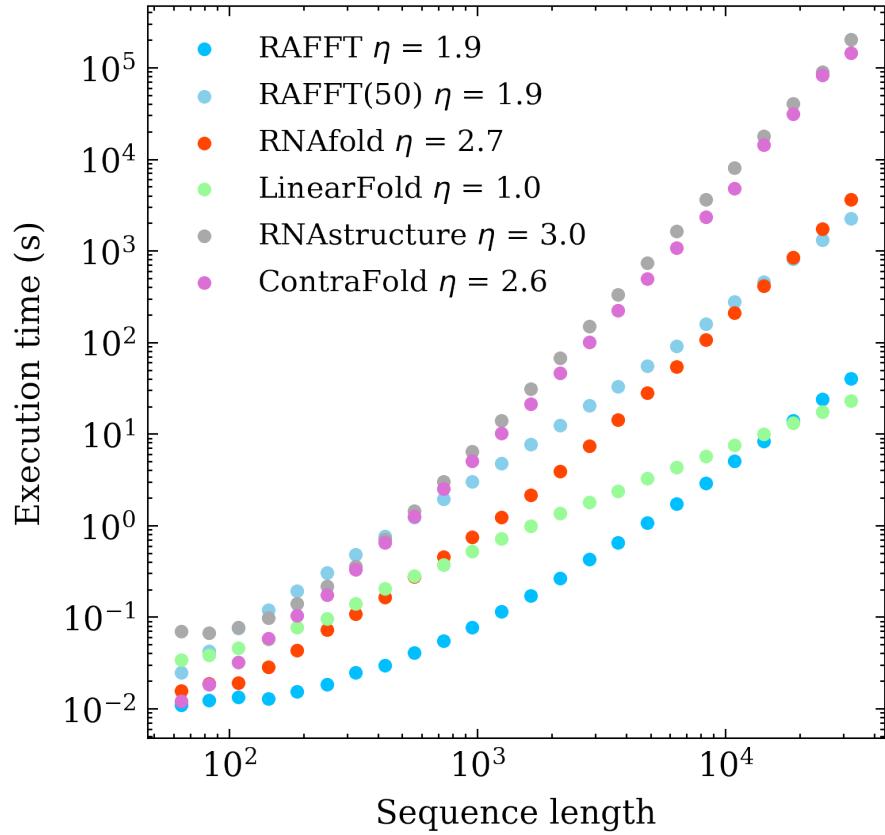
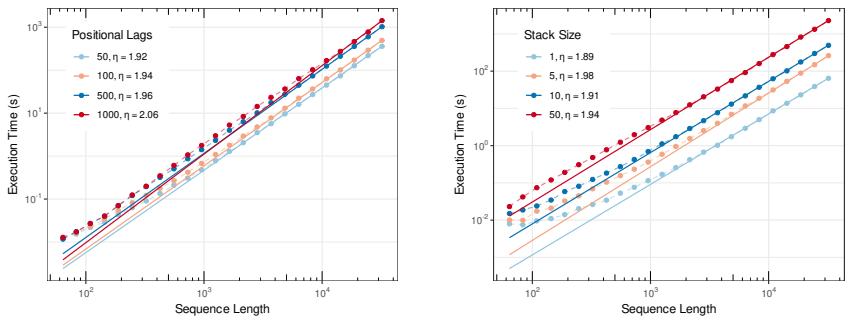


Figure 3.3: **Execution time comparisons.** For samples of 30 sequences per length, we averaged the execution times of five folding tools. The empirical time complexity $O(L^\eta)$ where η is obtained by non-linear regression. RAFFT denotes the naive algorithm (with only $N = 1$ structure saved per stack), whereas RAFFT(50) denotes the algorithm where 50 structures can be saved per stack.

2356 (exactly one base pair peer stems) can be formed. The approxi-
 2357 mate complexity therefore depends on

$$\sum_{i=0}^{L/2} (L - 2i) \log(L - 2i) = O(L^2 \log L) \quad (3.7)$$

2358 We compared RAFFT's execution time to the classical cubic-time
 2359 algorithms represented by CONTRAfold (Version 2.02), RNAstructure
 2360 (Version 2.0), RNAfold (Version 2.4.13) and the recent improved
 2361 DP tool LinearFold (Version 1.0). Figure 3.3 shows the execution
 2362 time of the RUST implementation of RAFFT and the four above-
 2363 mentioned tools for 30 random generated sequences of various
 2364 lengths. When comparing RAFFT implementation to the standard



(a) CPU times respect to the positional lags (n) (b) CPU times respect to the stack size (N)

Figure 3.4: Impact of the number of positional lags n and the stack size N on the runtime complexity. For a corresponding length, we generated 30 random sequences, and averaged their execution times. Solid lines indicate the estimated time complexity $O(L^\eta)$ where η is obtained with a non-linear regression on these average execution times for (A) Different number of positional lags n and (B) Different stack sizes N .

2365 DP tools, the execution time of RAFFT scales slower (with an ex-
 2366 ponent ≈ 2) with the sequence length whereas the standard DP
 2367 execution times are cubic. In contrast, the execution time of the
 2368 improved DP implemented by LinearFold scales linearly with
 2369 the sequence length. Only when considering a stack size of 1,
 2370 that RAFFT execution time is lower than the one of LinearFold
 2371 for sequence of lengths less than $L = 10^4$.

2372 We also analyse the scalability of RAFFT computational time
 2373 with respect to its critical parameters (the number of positional
 2374 lag n and the stack size N). Figure 3.4 shows for both different
 2375 stack sizes and number of positional lags, RAFFT execution time
 2376 against the sequence length. For both stack size and number of
 2377 positional lags, the execution time scales almost with the same
 2378 exponent (≈ 2).

2379 In sum, RAFFT’s performance shows a significant improvement
 2380 compared to three folding tools (RNAfold, RNAstructure, and
 2381 ContraFold), and we can approximate its theoretical time complex-
 2382 ity to $O(L^2 \log L)$, where L is the sequence length. However, its
 2383 average CPU time scales with respect to the stack size and the
 2384 number of positional lags considered. When $N = 1$ and $n = 100$,
 2385 RAFFT CPU time is lower than all of the four tools except for se-
 2386 quences longer than 10^4 . But when considering $N = 50$ stacks,
 2387 LinearFold showed better performance. Fitting the empirical
 2388 CPU times of each tool to a non-linear regression showed that all

2389 the methods scaled with respect to the sequence length whereas,
 2390 `LinearFold` scales linearly (i.e. $L = 1$) followed by RAFFT with an
 2391 exponent of $L \approx 2$, the MFE prediction methods scale cubically.
 2392 Now, does the improvement in CPU time guarantee the quality
 2393 of the predictions? The following section analyses the quality of
 2394 the structure predictions.

2395 *3.2.2 Accuracy of the predicted structural ensemble*

2396 After comparing RAFFT's computational time to existing tools, it
 2397 is also essential to assess the quality of the predicted secondary
 2398 structures. The quality of each tool's predictions is measured
 2399 using two statistical metrics: the PPV and the sensitivity. This
 2400 section presents the quality comparison of RAFFT predictions to
 2401 the four previously mentioned tools, i.e. RNAfold, RNAstructure,
 2402 LinearFold, Contrafold and the ML method Mfold2.

2403 We started by analyzing the prediction performances with re-
 2404 spect to sequence lengths: we averaged the performances at fixed
 2405 sequence length. Figure 3.5 shows the performance in PPV and
 2406 sensitivity for the five methods. It shows that the ML method
 2407 (Mfold2) consistently outperformed RAFFT and the other predic-
 2408 tions. When comparing only the MFE predictions produced using
 2409 the DP tools, LinearFold outperformed all other tools (RNAfold
 2410 and RNAstructure) for both short and long sequences. The *t*-test
 2411 between the ML and the most used MFE prediction tool (RNAfold)
 2412 revealed not only a significant difference (p -value $\approx 10^{-12}$) but
 2413 also a substantial improvement of 14.5% in PPV. RAFFT showed
 2414 performances similar to RNAfold; but, RAFFT is significantly less
 2415 accurate (p -value ≈ 0.0002), with a drastic loss of performance
 2416 for sequences of length greater than 300 nucleotides (See also
 2417 Table 3.1).

2418 However, are there relevant structures in the ensemble pre-
 2419 dicted by our method? To address this question we retained
 2420 the structure with the best score among the 50 recorded struc-
 2421 tures per sequence. We obtained an average PPV of 60.0% and
 2422 an average sensitivity of 62.8% over all the dataset. The gain in
 2423 terms of PPV/sensitivity is especially pronounced for sequences
 2424 of length ≤ 200 nucleotides, indicating the presence of biolog-
 2425 ically more relevant structures in the predicted ensemble than
 2426 the thermodynamically most stable one (PPV was =79.4%, and
 2427 sensitivity=81.2%). The average scores are shown in Table 3.1.
 2428 We also investigated the relation to the number of bases between

²⁴²⁹ paired bases (base pair spanning), but we found no striking effect,
²⁴³⁰ as already pointed out in one previous study [1].

Table 3.1: **Average performance displayed in terms of PPV and sensitivity.** The metrics were first averaged at fixed sequence length, limiting the over-representation of shorter sequences. The first two rows show the average performance for all the sequences for each method. The bottom two rows correspond to the performances for the sequences of length ≤ 200 nucleotides.

	RNAfold	LinearFold	RNAstructure	CONTRAfold	Mxfold2	RAFFT	RAFFT*
All sequences							
PPV	55.9	60.6	54.7	58.4	70.4	47.7	60.0
Sensitivity	63.3	58.9	61.5	65.2	77.1	52.8	62.8
Sequences with lengths ≤ 200							
PPV	59.5	63.2	58.2	60.5	76.7	57.9	79.4
Sensitivity	65.5	59.4	63.8	65.9	82.9	63.2	81.2

²⁴³¹ All methods performed poorly on two groups of sequences:
²⁴³² one group of 80 nucleotides long RNAs, and the second group of
²⁴³³ around 200 nucleotides (three examples of such sequences are
²⁴³⁴ shown in the Appendix A3.1). Both groups have large unpaired
²⁴³⁵ regions, which for the first group lead to structures with average
²⁴³⁶ free energies 9.8 kcal/mol according to our dataset. The PCA anal-
²⁴³⁷ ysis of the native structure space, shown in Figure 3.6, reveals a
²⁴³⁸ propensity for interior loops and the presence of large unpaired
²⁴³⁹ regions like hairpins or external loops. Figure 3.6 shows the struc-
²⁴⁴⁰ ture space produced by Mxfold2, which seems close to the native
²⁴⁴¹ structure space. In contrast, the structure spaces produced by
²⁴⁴² RAFFT and RNAfold are similar and more diverse.

²⁴⁴³ In summary, we performed the prediction quality compari-
²⁴⁴⁴ son for different sequence lengths. The dataset was divided into
²⁴⁴⁵ two sets: one with lengths less than 200 nucleotides and the rest
²⁴⁴⁶ constituting the second. Because RAFFT predicts an ensemble of
²⁴⁴⁷ structures, which contrasts the other tools, we also distinguish
²⁴⁴⁸ the single prediction (RAFFT) comparison from the ensemble
²⁴⁴⁹ one (RAFFT*). Overall, on average, RAFFT performed qualitatively
²⁴⁵⁰ poorer than existing tools in terms of both PPV and sensitivity. The
²⁴⁵¹ ML method, Mxfold2 outperformed all existing methods for differ-
²⁴⁵² ent RNA sequence lengths but equalized RAFFT* performance for
²⁴⁵³ sequences of length less than 200 nucleotides. The later showed

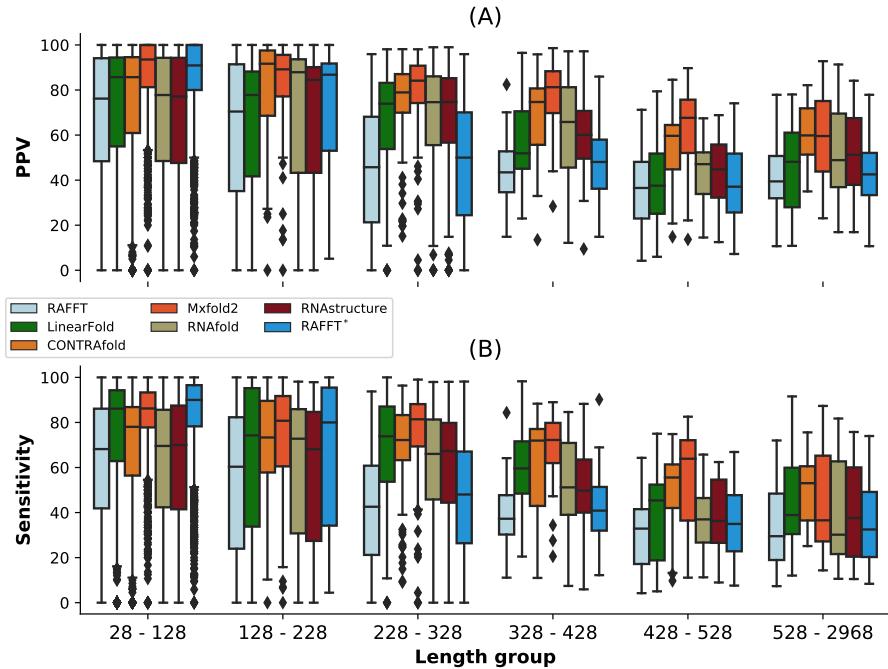


Figure 3.5: **RAFFT’s performance on folding task.** (A) PPV vs sequence length. In the top panel, RAFFT (in light blue) shows the PPV score distributions when for the structure (out of $N = 50$ predictions) with the lowest free energy, whereas RAFFT* (in blue) shows the best PPV score in that ensemble. (B) Sensitivity vs sequence length.

that RAFFT predicted ensemble contains sequences of biological interest. We further assess the quality of that ensemble with the proposed kinetics ansatz. The next section discusses two RNA kinetic test cases: the application of the kinetic ansatz on CFSE and a bistable RNA sequence.

3.2.3 Applications to the RNA kinetics

Furthermore, the ensemble of structures predicted by RAFFT is analyzed using a kinetics ansatz to extract information about the dynamic of RNA folding. This section analyses the kinetics of two RNA sequences using RAFFT predicted pathways.

We started with the CFSE, a natural RNA sequence of 82 nucleotides with a structure determined by sequence analysis and obtained from the RFAM database. This structure has a pseudo-knot which is not taken into account here.

[Figure 3.7](#)A and [Figure 3.7](#)B show respectively the fast-folding graph constructed using RAFFT, and the MFE and native structures

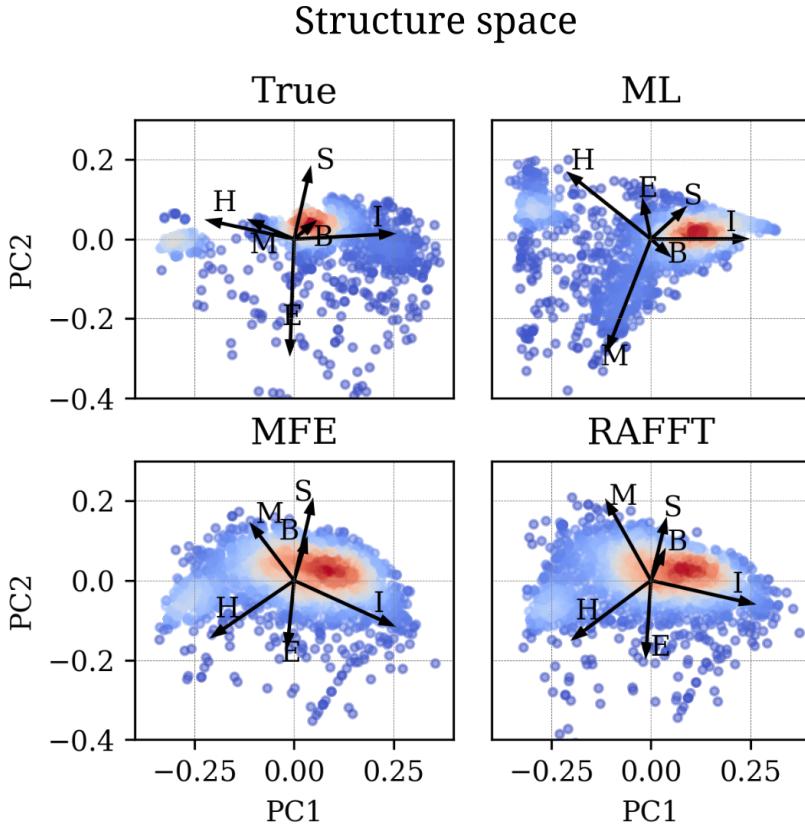


Figure 3.6: **Structure space analysis.** PCA for the predicted structures using RAFFT, RNAfold, MxFold2 compared to the known structures denoted “True”.

for the CFSE. The fast-folding graph is computed in four steps. At each step, stems are constructed by searching for $n = 100$ positional lags and, a set of $N = 20$ structures (selected according to their free energies) are stored in a stack. The resulting fast-folding graph consists of 68 distinct structures, each of which is labelled by a number. Among the structures in the graph, 6 were found similar to the native structure (16/19 base pairs differences). The structure labelled “29” in the graph leading to the MFE structure “59” is the 9th in the second stack. When storing less than 9 structures in the stack at each step, we cannot obtain the MFE structure using RAFFT; this is a direct consequence of the greediness of the proposed method. To visualize the energy landscape drawn by RAFFT, we arranged the structures in the fast-folding graph onto a surface according to their base-pair distances; for this we used the multidimensional scaling algorithm implemented in the `scipy` package. Figure 3.7D shows the landscape interpolated with all the structures found; this landscape illustrates the bi-stability

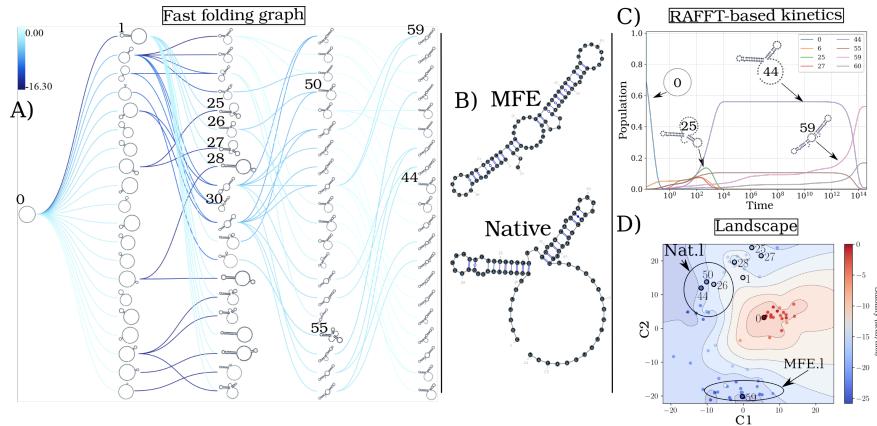


Figure 3.7: Application of the folding kinetic ansatz on CFSE. (A) Fast-folding graph in four steps and $N = 20$ structures stored in a stack at each step. The edges are coloured according to $\Delta\Delta G$. At each step, the structures are ordered by their free energy from top to bottom. The minimum free energy structure found is at the top left of the graph. A unique ID annotates visited structures in the kinetics. For example, “59” is the ID of the **MFE** structure. (B) **MFE** (computed with RNAfold) and the native **CFSE** structure. (C) The change in structure frequencies over time. The simulation starts with the whole population in the open-chain or unfolded structure (ID 0). The native structure (**Nat.1**) is trapped for a long time before the **MFE** structure (**MFE.1**) takes over the population. (D) Folding landscape derived from the 68 distinct structures predicted using RAFFT. The axes are the components optimized by the MDS algorithm, so the base pair distances are mostly preserved. Observed structures are also annotated using the unique ID. **MFE**-like structures (**MFE.1**) are at the bottom of the figure, while native-like (**Nat.1**) are at the top.

2487 of the **CFSE**, where the native and **MFE** structures are in distinct
 2488 regions of the structure space.

2489 From the fast-folding graph produced using RAFFT, the transi-
 2490 tion rates from one structure in the graph to another are computed
 2491 using the formula given in [Equation 3.6](#). Starting from a popu-
 2492 lation of unfolded structure and using the computed transition
 2493 rates, the native of structures is calculated using [Equation 2.3](#).
 2494 [Figure 3.7C](#) shows the frequency of each structure; as the fre-
 2495 quency of the unfolded structure decreases to 0, the frequency of
 2496 other structures increases. Gradually, the structure labelled “44”,
 2497 which represents the **CFSE** native structure, takes over the popu-
 2498 lation and gets trapped for a long time, before the **MFE** structure
 2499 (labelled “59”) eventually becomes dominant. Even though the

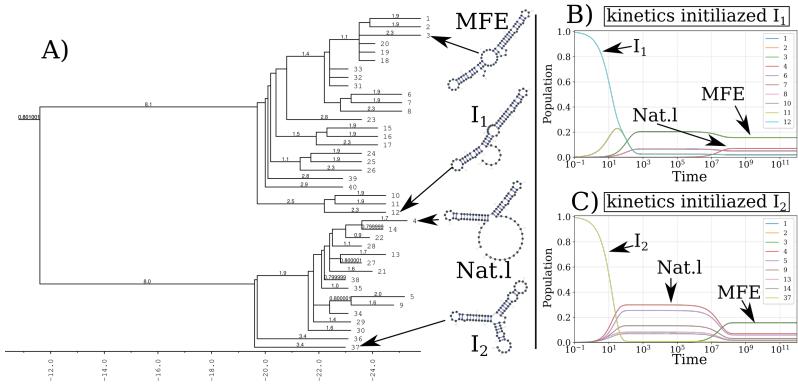


Figure 3.8: Folding kinetics of CFSE using Treekin. A) Barrier tree of the CFSE. From a set of 1.5×10^6 sub-optimal structures, 40 local minima were found, connected through saddle points. The tree shows two alternative structures separated by a high barrier with the global minimum (**MFE** structure) on the right side. (B) Folding kinetics with initial population I_1 . Starting from an initial population of I_1 , as the initial frequency decreases, the others increase, and gradually the **MFE** structure is the only one populated. (C) Folding kinetics with initial population I_2 . When starting with a population of I_2 , the native structure (labelled **Nat.1**) is observable, and gets kinetically trapped for a long time due to the high energy barrier separating it from the **MFE** structure.

fast-folding graph does not allow computing energy landscape properties (saddle, basin, etc.), the kinetics built on it reveals a high barrier separating the two meta-stable structures (**MFE** and native).

Our kinetic simulation was then compared to Treekin [55]. First, we generated 1.5×10^6 sub-optimal structures up to 15 kcal/mol above the **MFE** structure using RNAsubopt [112]. Since the **MFE** is $\Delta G_s = -25.8$ kcal/mol, the unfolded structure could not be sampled. Second, the ensemble of structures is coarse-grained into 40 competing basins using the tool barriers [55], with the connectivity between basins represented as a barrier tree (see Figure 3.8A). When using Treekin, the choice of the initial population is not straightforward. Therefore we resorted to two initial structures I_1 and I_2 (see Figure 3.8B and 3.8C, respectively). In Figure 3.8B, the trajectories show that only the kinetics initialized in the structure I_2 can capture the complete folding dynamics of CFSE, in which the two metastable structures are visible. Thus, in order to produce a folding kinetics in which the native and the **MFE** structures are visible, the kinetic simulation performed using

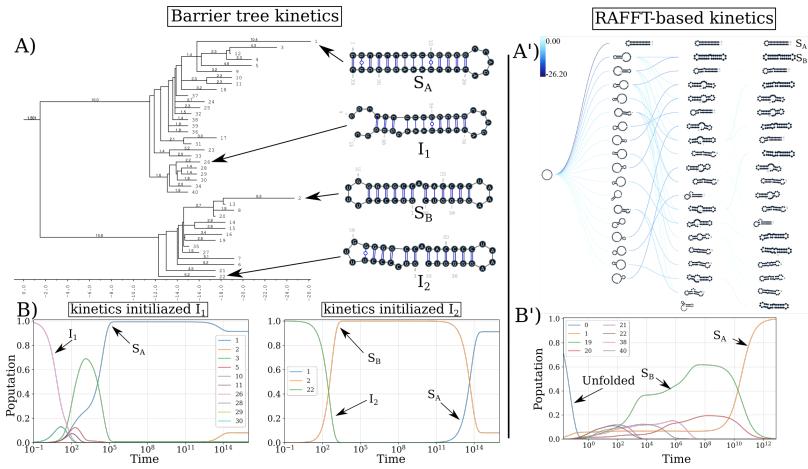


Figure 3.9: RAFFT vs Treekin: folding kinetics of a bi-stable RNA sequence. (A) Barrier tree for the bi-stable example sequence. The local minima and the corresponding barriers are computed from the complete enumeration of the structure space. The bi-stability is visible on the barrier tree through the two branches separated by a high barrier. (B) Folding kinetics trajectories. The left plot shows the folding dynamics starting from a population with I_1 , and the right size is the kinetics when the population is initialized in structure I_2 . When starting from I_1 , S_A is quickly populated; starting from I_2 , the bi-stability is more apparent. (A') Fast-folding graph using RAFFT. A maximum of $N = 20$ structures are stored in a stack at each step and overall 46 distinct structures are visited. (B') Folding kinetics trajectory obtained from the fast-folding graph (indices are different from the barrier tree indices). The dynamics starts with a population with only unfolded structure, and slowly, S_B is populated and gets trapped for a long time before the MFE structure S_A becomes populated.

2519 Treekin required a particular initial condition and a barrier tree
 2520 representation of the energy landscape built from a set of 1.5×10^6
 2521 structures. By contrast, using the fast-folding graph produced by
 2522 RAFFT, which consists only of 68 distinct structures, our kinetic
 2523 simulation produces complete folding dynamics starting from a
 2524 population of unfolded structure.

2525 As a second illustrative example, we applied both kinetic mod-
 2526 els to the classic bi-stable sequence. For Treekin, we first sam-
 2527 pled the whole space of 20×10^3 sub-optimal structures from
 2528 the unfolded state to the MFE structure, and from that set, 40
 2529 basins were also computed using barriers. The barrier tree in
 2530 Figure 3.9 shows the bi-stable landscape, where the two deepest

2531 minima are denoted S_A and S_B . As in the first application, we also
 2532 chose two initializations with the structures denoted I_1 and I_2 in
 2533 [Figure 3.9A](#) and [3.9B](#). Secondly, we simulate the kinetics starting
 2534 from the two initial conditions (See [Figure 3.9B](#)). When starting
 2535 from I_2 , the slow-folding dynamics is visible: S_B first gets kineti-
 2536 cally trapped, and the [MFE](#) structure (S_A) only takes over later on.
 2537 For our kinetic ansatz, we started by constructing the fast-folding
 2538 graph using RAFFT, consisting of only 46 distinct structures. The
 2539 resulting kinetics, shown in [Figure 3.9B'](#) was found qualitatively
 2540 close to the barrier kinetics initialized with structure I_2 . Once
 2541 again, with few as 48 structures, our proposed kinetic ansatz can
 2542 produce complete folding dynamics starting from a population
 2543 of unfolded structure.

2544 In both examples, our kinetic ansatz derived from the fast fold-
 2545 ing graph predicted by RAFFT produces complete folding kinetic
 2546 trajectories, using fewer structures than the existing methods that
 2547 required the complete enumeration of the fitness landscape (i.e.
 2548 all structures and their associated energies). Despite the poor
 2549 validation procedure of our kinetic ansatz, we believe that the
 2550 [RNA](#) pathways predicted by RAFFT could contain structures of bio-
 2551 logical pertinence. An analysis of the sample structures produced
 2552 by RAFFT is provided in Appendix [Section A.1](#) and a discussion
 2553 on some limitations in [Chapter 6](#).

2554 3.3 CONCLUSION

2555 We have proposed a method for [RNA](#) structure, and dynamics
 2556 predictions called RAFFT. Our method is inspired by the experi-
 2557 mental observation of parallel fast-folding pathways. To mimic
 2558 this observation, we designed an algorithm that produces par-
 2559 allel folding pathways in which stems are formed sequentially.
 2560 Taking advantage of the [FFT](#), the time complexity of our method
 2561 was slowed down to $O(L^2 \log L)$, thus improving the cubic time
 2562 complexity of classic [DP](#) methods. Then, we proposed a kinetic
 2563 ansatz that exploits the parallel fast-folding pathways predicted
 2564 to model how different conformations are populated over time.
 2565 Our kinetic ansatz produced complete folding dynamics without
 2566 sampling the entire conformation space. However, our method
 2567 also presents some limitations that will be discussed in [Chapter 7](#).

[August 4, 2022 at 19:02 – 1.0]

2568

Part II

2569

RNA DESIGN

2570

This second part of our thesis fucuses only on the inverse folding of RNA secondary structures. It contains figures and ideas that have previously appeared in our publications:

2574

- [128] **Nono SC Merleau** and Matteo Smerlak (2021). *A simple evolutionary algorithm guided by local mutations for an efficient RNA design*. In: Proceedings of the Genetic and Evolutionary Computation Conference. pp. 1027-1034. (Published)
- [129] **Nono SC Merleau** and Matteo Smerlak (2022). *An evolutionary algorithm for inverse RNA folding inspired by Lévy flights* In: *bioRxiv* (Submitted and accepted) (BMC Bioinformatics).

[August 4, 2022 at 19:02 – 1.0]

2583

2584 INTRODUCTION TO RNA DESIGN

2585 The previous chapters demonstrated the implications of ncRNAs
2586 molecules in varying levels of cellular processes, from gene ex-
2587 pression regulation (miRNAs, piRNAs, lncRNAs) to RNA matura-
2588 tion (sncRNAs, snoRNAs) and protein synthesis (rRNAs, tRNAs).
2589 Knowing that these biological functions are performed by high
2590 dimensional RNA structures, which strongly depend on their sec-
2591 ondary structures, we also provided a comprehensive review
2592 of computation methods for predicting secondary structures.
2593 Now that we have computational folding tools that are accu-
2594 rate enough, is it possible to design an RNA molecule that can
2595 accomplish a desired biological function for a given secondary
2596 structure? Answering this question may demand both experi-
2597 mental and computational efforts. For artificial ncRNAs for which
2598 the native RNA sequence is unknown, the essential prerequisite
2599 for experimentalists is often a computational solution to the in-
2600 verse folding problem. Unlike the folding situation, the inverse
2601 folding problem begins from a given secondary structure, and
2602 the goal is to find one or many RNA sequences that fold into that
2603 secondary structure. This chapter aims to provide the formal
2604 background and biotechnological implications of addressing this
2605 problem. Then, it gives a brief literature review of the existing
2606 computational methods.

2607 4.1 RNA INVERSE FOLDING AND BIOTECHNOLOGICAL IMPLICA-
2608 TIONS

2609 In modern biotechnology, we often seek to reproduce the natural
2610 ability of the cells to control gene expressions using a variety of
2611 nucleic acids and proteins. These natural cellular abilities result
2612 from networks of regulatory molecules such as ncRNAs that dy-
2613 namically regulate the expression of specific genes in response to
2614 environmental signals. Therefore, the ability to engineer biologi-
2615 cal systems is directly related to controlling gene expression. The
2616 increasing number of examples of natural regulator ncRNAs has
2617 opened doors to many emerging subfields such as RNA synthetic
2618 biology [21, 86] and RNA nanostructure [68, 89]. Researchers
2619 have engineered RNA molecules with new biological functions, in-

2620 spired by this natural versatility. Synthetic biology has also made
 2621 significant progress in developing versatile and programmable
 2622 genetic regulators that precisely control gene expressions in the
 2623 last decades. Three general approaches are taken to engineer new
 2624 functional RNAs: harvesting from nature, computational design
 2625 and molecular evolution. We are interested here in computational
 2626 RNA design methods.

2627 In most cases, designing a functional RNA goes beyond com-
 2628 putationally generating a set of RNA sequences that fold into a
 2629 given secondary structure. Successful design methods include
 2630 computational and experimental, predictive and analytical tech-
 2631 niques. However, computational tools addressing the inverse
 2632 folding problem often provide some guidance and rationalities
 2633 through the design process. For example, Steffen Mueller and
 2634 his collaborators [134] suggested a systematic, rational approach,
 2635 synthetic attenuated virus engineering (SAVE), to develop new,
 2636 productive live attenuated influenza virus vaccine candidates
 2637 using computer-aided rational design. In addition, Eckart Binde-
 2638 wald et al. [13] used computational tools for solving inverse RNA
 2639 folding in the design of nanostructures, including pseudoknots.
 2640 And in designing several ncRNAs with a successful synthetic such
 2641 as ribozymes [40], riboswitches [52, 203].

2642 Depending on the specificities of the RNA design task, finding
 2643 the underlying mathematical model that maps each designed
 2644 RNA sequence solution to a set of properties that includes most of
 2645 the specifications or constraints can be a challenging task. When it
 2646 exists, it allows to address the RNA design problem computa-
 2647 tionally, and we call this mathematical model the objective function of
 2648 the RNA design problem. The complexity of the objective function
 2649 used gives rise to two RNA design problems: the negative and the
 2650 positive design. The following section describes both RNA design
 2651 problems and their computational complexities.

2652 4.2 THE POSITIVE AND NEGATIVE DESIGN

2653 We often find two types of RNA design problems in the literature:
 2654 negative and positive design. The negative structural design of
 2655 RNAs, also called the inverse RNA folding problem, aims to find
 2656 one or many RNA sequences that fold into a given target RNA
 2657 secondary structure while avoiding alternative folds of similar
 2658 quality for the chosen energy model ΔG . In other terms, it is an
 2659 optimization problem where a target RNA secondary structure S^*

2660 of length L is given, and the goal is to determine an [RNA](#) sequence
 2661 ϕ of length L such that $\forall \mathcal{S} \neq \mathcal{S}^* \in \Sigma_\phi, \Delta G(\phi, \mathcal{S}) > \Delta G(\phi, \mathcal{S}^*)$.

2662 This problem is [NP-hard](#) even in a simple energy model [14],
 2663 and we cannot provide a parameterized algorithm that solves it
 2664 in a polynomial time.

2665 In contrast, a positive design problem consists of optimizing
 2666 affinity towards a given target secondary structure. In other terms,
 2667 the objective is to find a sequence $\phi \in \{A, U, C, G\}^L$ such that
 2668 $\mathcal{S}^* = \mathcal{S}^{MFE}(\phi) = \arg \min_{\mathcal{S} \in \Sigma_\phi} \Delta G(\phi, \mathcal{S})$ (i.e. the sequence ϕ
 2669 should have as [MFE](#)s structure of its ensemble Σ_ϕ the target struc-
 2670 ture \mathcal{S}^*). The positive design is computationally solvable exactly
 2671 in polynomial time [54].

2672 Both negative and positive designs are considered in this work,
 2673 and the main difference often depends on the objective function
 2674 used. In addition, it has been recently shown that the proportion
 2675 of designable secondary structures decreases exponentially with
 2676 L for various popular combinations of energy models and design
 2677 objectives [216]. The following section presents an overview of
 2678 previously used objective functions for the [RNA](#) design prob-
 2679 lem.

2680 4.3 OBJECTIVE FUNCTIONS PREVIOUSLY USED IN THE CONTEXT 2681 OF INVERSE RNA FOLDING

2682 For a given target secondary structure \mathcal{S}^* of length L , a brute force
 2683 approach to the inverse [RNA](#) folding problem that enumerates
 2684 all possible [RNA](#) sequences is not viable due to the exponential
 2685 growth of the search space with increasing length (i.e. 4^L). For
 2686 the space of compatibles sequences to the target \mathcal{S}^* , an upper
 2687 bound can be defined by restricting the paired position to the
 2688 base pairs: G-C, G-U, and A-U. This results in $6^{(L-u)/2} \times 4^u$ se-
 2689 quences compatible with \mathcal{S}^* where u is the number of unpaired
 2690 nucleotides. The most common way to efficiently handle the huge
 2691 set of possible solutions is to solve an optimization problem sub-
 2692 jected to a formulated objective function. There exists a variety of
 2693 well-established optimization methods helping to perform this
 2694 task. However, finding the right objective function to evaluate
 2695 the solutions can be quite challenging. This section of our work
 2696 provides an overview of an objective function and an essential
 2697 description of the most previously used objective functions in
 2698 designing [RNA](#) molecules.

2699 The objective function defines a mathematical model that maps
 2700 each [RNA](#) sequence solution to its essential properties or func-

tions. In biological terms, this relation between fitness and sequence can be seen as assigning a phenotype (score) to a genotype (sequence). Selection pressure due to the optimization method ensures that better phenotypes are advantageous and thus preferred, which optimizes the sequence to fall into fitness optima. This section defines the previously used objective functions in the [RNA](#) design problems and highlights some interesting properties.

- A simple distance from the target structure: in the simplest setting, the objective function of an [RNA](#) sequence ϕ defines the distance between \mathcal{S}^* and the current [MFE](#) structure $\mathcal{S}^{MFE}(\phi)$. It often requires only the [MFE](#) structure's computation, hence being computationally fast. There are many variants of this distance measure: base-pair distance, hamming or string edit distance, tree-edit distance and energy distance. For a formal definition of each of those distances, see [Section 1.4](#). This objective function was used in the earliest tools such as [RNAINVERSE](#) [80] but also in many others since then [5, 18, 59].
- A negative design objective function: in contrast to the above mentioned objective functions (often considered when performing a positive design), we consider the whole structural ensemble when computing the fitness of an [RNA](#) sequence ϕ . In most cases, it is preferable also to consider negative design goals, which allows for avoiding alternative structures of similar quality to the target structure. Negative [RNA](#) design methods usually consider one of the three following defects: (1) the *suboptimal defect* [37, 54, 80, 218] which defines the energy distance to the first suboptimal (2) the *probability defect* [80, 218] which defines the probability that the sequence ϕ folds into any other structure than the target structure \mathcal{S}^* and (3) the *ensemble defect* [218] which corresponds to the average number of incorrectly paired nucleotides at equilibrium calculated over the structure ensemble of ϕ , Σ_ϕ .
- Multi-objective optimization: in some designing cases where more than one goal is specified, it is necessary to formulate an objective function for each goal. That results in a multi-objective optimization problem. The solutions to such a problem are all optimal for at least one objective function and thus arranged on the so-called Pareto optimal front. This approach has already been used in several [RNA](#) design tools such as [Modena](#) [189, 190] and in [143].

- Bistable and multi-stable riboswitches objective functions:
In some designing cases, especially for riboswitches, it is possible to specify more than one desired target structure, including the energy differences between them, the barrier heights and the kinetic properties. Following the same idea, Flamm et al. introduced an objective function that enables designing RNA molecules to adopt two distinct structures [54]. This bistable objective function contains two terms. The first term increases the probability of both structures in the ensemble, and the second specifies the desired energy difference between both states. It is also possible to vary the states' temperature to gain a bistable thermoswitch. The same idea has therefore been expanded to an objective function for designing RNA molecules that can adopt more than two structures, including extension for multi-structure energy barrier calculations [143, 173]. Frnakenstein [117] also utilises such objective function for multi-target design.
- Mutational robustness and neutrality: In addition to the above-mentioned objective functions, objective functions aim to measure the mutual neutrality of the sequence concerning the target structure [173]. When using such an objective function, the sequences are optimized so that the fraction of one-mutant neighbours to the original structure is as significant as possible. This allows for perfectly preserving the structure when mutations are introduced. We often talk of a mutational robustness optimization [6].

These objective functions suggest that the inverse folding problem is a major challenge with no single solution yet, and many possible ways of setting the goal. This thesis relies on three objective functions: the simple distance to the target, the ensemble defect, and the mutational robustness. In addition to the many objective functions, there are also several methods. The following section will review the existing methods independently of the objective function and provide some limitations.

2777 4.4 A REVIEW ON EXISTING INVERSE RNA FOLDING TOOLS.

2778 Several methods or algorithms addressing this problem have
2779 been proposed in the literature. The existing techniques can be
2780 classified into two categories: one for the pseudoknot-free struc-
2781 ture design and another for the pseudoknotted RNA structure

2782 design. This section gives a short description of some of the exist-
 2783 ing tools, especially those used in the benchmark results of the
 2784 thesis.

2785 4.4.1 *Pseudoknot-free RNA inverse folding tools*

2786 Due to the complexity of the [RNA](#) design, most of the existing
 2787 tools perform a stochastic search optimization where initial po-
 2788 tential solutions are generated and refined over a finite number
 2789 of iterations or generations [43, 47, 48, 141, 188]. Some stochas-
 2790 tic search techniques may involve several candidate solutions at
 2791 each generation or not. The ones that do are population-based
 2792 algorithms, which means they maintain a set of candidate solu-
 2793 tions at each generation, with each solution corresponding to a
 2794 unique point in the problem’s search space. We are interested
 2795 in this work in [EA](#), a particular class of population-based algo-
 2796 rithms. This section presents an overview of [EA](#) when applied to
 2797 the inverse folding of [RNA](#) molecules. In addition, it reviews the
 2798 existing tools implementing similar and different techniques.

2799 4.4.1.1 *Evolutionary algorithms and [RNA](#) inverse problems*

2800 Among the existing tools dealing with the [RNA](#) inverse problem,
 2801 both [ERD](#) [47, 48] and [MODENA](#) [189] are [EAs](#) but implementing
 2802 different strategies. In general, an evolutionary search algorithm
 2803 on any fitness landscape consists of three main parts, which in
 2804 the context of [RNA](#) inverse folding are as follows:

- 2805 • Initialization: generating a random initial population of
 2806 [RNA](#) sequences compatible with the given target secondary
 2807 structure.
- 2808 • Evaluation and selection: evaluating a population of [RNA](#)
 2809 sequences consists of two steps: 1) fold each sequence into
 2810 a secondary structure and assign it a weight based on its
 2811 similarity to the target structure. 2) select a weighted ran-
 2812 dom sample with replacement from the current population
 2813 to generate a new population. A detailed description of
 2814 the objective function used in our proposed tool [aRNAque](#) is
 2815 provided in the next chapter.
- 2816 • Mutation (or move) operation: define a set of rules or steps
 2817 used to produce new sequences from the selected or ini-
 2818 tial ones. This component is elaborated further in the next
 2819 chapter.

2820 MODENA uses a multi-objective function that measures the stabili-
 2821 ty of the folded sequence and its similarities to the target. It starts
 2822 from a population of randomly generated sequences, and the ob-
 2823 jective is optimized through tournament selection and random
 2824 mutation at non-closing loop positions.

2825 In contrast, ERD starts by decomposing the target structure
 2826 into loops and independently uses an evolutionary algorithm
 2827 to minimize each constituent's energy. It was first developed in
 2828 2014 [48], and one year after, an updated version was released
 2829 [47]. The main lines of ERD are:

- 2830 1. Pool reconstruction: using a collection of **RNA** sequences
 2831 (STRAN database) similar to the natural ones, a pool of
 2832 sequences is constructed for their length by successively
 2833 finding the corresponding structure using **RNAfold**, decom-
 2834 posing the structure in sub-components, and finally, the
 2835 corresponding sub-sequences of the same size are gathered
 2836 to form a pool.
- 2837 2. Hierarchical decomposition of the target structure into loops:
 2838 using the idea that any secondary structure can be uniquely
 2839 decomposed into its structural components (stems, hairpin
 2840 loops, internal loops, bulge and multi-loops), ERD decom-
 2841 poses the target in the positions where multi-loops occur.
- 2842 3. Sequence initialization: after decomposing the target struc-
 2843 ture into sub-components, for each sub-component, a ran-
 2844 dom sub-sequence is chosen from the pool, and the initial
 2845 sequence is a combination of those sub-sequences;
- 2846 4. Evolutionary optimization of the sub-sequences: an EA al-
 2847 gorithm is performed on each sub-component to improve
 2848 the initial sequence. The outcome sub-sequences are com-
 2849 bined to form a newer sequence that will replace the initial
 2850 one. Iteratively the evolutionary algorithm is performed on
 2851 the updated sequence until the combined sequence folds
 2852 into the target or in a failure case when the stopping con-
 2853 dition is satisfied. Two evolutionary operators are imple-
 2854 mented here, a mutation that consists of replacing a sub-
 2855 sequence corresponding to a sub-component with a new
 2856 random one from the pool for the same length, and a se-
 2857 lection which consists of choosing from a population of 15
 2858 **RNA** sequences or sub-sequences, three best sequences with
 2859 respect to their free energy and adding them to the best
 2860 from the previous generation, three best ones with respect

2861 to the Hamming distance from the target are therefore cho-
 2862 sen. The next-generation population is then obtained by
 2863 generating five new sequences for each of the three best
 2864 sequences.

2865 In the different EA methods presented above, the mutation
 2866 operation is essential for good performance because it provides
 2867 the rules that allow for navigating the solution space. ERD imple-
 2868 ments a non-local mutation, which consists of randomly changing
 2869 a subsequence in the candidate solution with a new one taken
 2870 from a set of possible moves. In contrast, Modena uses both local
 2871 mutation and crossover operation to improve its search. However,
 2872 both EAs present difficulties in finding RNA sequences that fold
 2873 into some secondary structures of the Eterna100 data set. That
 2874 limitation could be due to the local search (for Modena) or the
 2875 finite set of move data used in the non-local search implemented
 2876 in ERD. In mathematical optimization, local searches are known
 2877 for their quick convergence to a local minimum. This could be
 2878 the same case for EAs implementing local mutations. To avoid
 2879 early convergence EA practitioners often implement non-local
 2880 mutation methods, e.g. Lévy search, inspired by the Lévy flights.
 2881 The following section describes the Lévy flight and reviews some
 2882 applications of Levy search in the context of EAs.

2883 4.4.1.2 Lévy flights and evolutionary algorithms

2884 In this section, we define concepts such as Lévy flights and pro-
 2885 vide a brief review of its implications and applications to opti-
 2886 mization techniques such as evolutionary algorithms.

2887 In its classical setting, evolutionary algorithms are guided by
 2888 local (or one-point mutations) mutations. Although a local search
 2889 can efficiently discover optima in a simple landscape, more com-
 2890 plex landscapes pose challenges to designing evolutionary al-
 2891 gorithms that rely solely on local search. This is especially true
 2892 on a landscape with high neutrality where local search may be
 2893 inefficient or risk getting stuck on a plateau (or local optimum).
 2894 To avoid this pitfall, many practitioners suggested EA that imple-
 2895 ments a mutation scheme inspired by Lévy flights (called Lévy
 2896 mutation).

2897 Lévy flights are random walks with a Lévy (or any heavy-
 2898 tailed) step size distribution. The concept originates in the work of
 2899 Mandelbrot on the fluctuation of commodities prices in the 1960s
 2900 [120] but has since found many more physical applications [172].
 2901 The term "Lévy flight" was also coined by Mandelbrot, who

2902 used one specific distribution of step sizes (the Lévy distribution,
 2903 named after the French mathematician Paul Lévy). Lévy flights
 2904 also play a key role in animal foraging, perhaps because they
 2905 provide an optimal balance between exploration and exploitation
 2906 [93, 201]. For a recent review of applications of Lévy flights in
 2907 biology from the molecular to the ecological scale, [149].

2908 Similar to a Lévy flight, a Lévy mutation scheme allows simu-
 2909 taneous search at all scales over the landscape. New mutations
 2910 most often produce nearby sequences (one-point mutations), but
 2911 occasionally generate mutant sequences which are far away in
 2912 genotype space (macro-mutations). In this work, the distribution
 2913 of the number of point mutations at every step is taken to follow
 2914 a Zipf distribution [135].

2915 Earlier works have applied similar ideas in genetic program-
 2916 ming [30], and in differential evolutionary algorithms [170]. This
 2917 motivated us to investigate a possible benefit of a Lévy flight
 2918 in the design of RNA sequences in Chapter 5. In addition to EA
 2919 methods, there exists several computational RNA design tools
 2920 implementing different techniques such as, ML, nested monte
 2921 carlo search (NMCS) etc... The following section provides a short
 2922 description of such tools.

2923 4.4.1.3 Tools implementing non-EA strategies.

2924 Several tools dealing with the RNA folding problem implement
 2925 different strategies from the population-based, or evolutionary
 2926 algorithm approaches. This section describes couple of them,
 2927 emphasising on those that are used in the benchmark results in
 2928 Chapter 5, which are NEMO, RNAinverse, antaRNA and sentRNA.

2929 sentRNA [171] is a computational agent that uses a set of infor-
 2930 mation and strategies collected from the EteRNA game players
 2931 to train a neural network model. The neural network assigns an
 2932 identity of A, U, C, or G to each position in the given target, a
 2933 featured representation of its local environment. The featured
 2934 representation combines information about its bonding partner,
 2935 nearest neighbours, and long-range features. While the bonding
 2936 partner and nearest neighbour information are provided to the
 2937 agent by default, long-range features are learned through the
 2938 training data. For each target structure, the long-range features
 2939 refer to the important position j relative to i that the agent should
 2940 know about when deciding what nucleotide to assign to i . These
 2941 are defined by two values: the Cartesian distance and the angle in
 2942 radians. Those two values are computed for each position (i, j) us-

2943 ing a mutual information metric over the player solution dataset.
 2944 Therefore, the result is a list of long-range features for a given
 2945 target structure. A subset of long features is selected from this
 2946 list and used to define a model for the neural network model's
 2947 training, validation, and testing. In addition to the neural net-
 2948 work model, sentRNA also implements a refinement algorithm
 2949 on the unsuccessful design. The refinement algorithm is an adap-
 2950 tive walk that starts from the predicted sequence and uses a set
 2951 of random mutations that allow improving the neural network
 2952 solution. Alternatively, EternaBrain [105] implement a convo-
 2953 lutional network model trained on a huge EteRNA moves-select
 2954 repository of 30,477 moves from the top 72 players; and LeaRNA
 2955 [155] uses deep reinforcement learning to train a policy network
 2956 to sequentially design an entire RNA sequence given a specified
 2957 target structure.

2958 NEMO [141] is a recently developed tool combining a NMCS tech-
 2959 nique with domain-specific knowledge to create a novel algo-
 2960 rithm. The underlying idea is to start with an input pattern se-
 2961 quence of N's of the same length as the targeted structure. First,
 2962 it uses the standard NMCSs to sample sequence solutions acting
 2963 on N's only. A sequence candidate is selected from the sample;
 2964 then folded into an MFE structure. When the MFE structure does
 2965 not match the target, some subset mutations are performed, and
 2966 a set of random mutated positions are picked to generate a new
 2967 input pattern sequence. The new input pattern will allow sam-
 2968 pling acting on N's only using the same standard NMCSs. This
 2969 procedure is then repeated several times until the MFE structure
 2970 matches the targeted structure or not in the unsuccessful cases.
 2971 The statistical results show that NEMO surpasses all the existing
 2972 tools on the EteRNA100 benchmark datasets by solving $\approx 95\%$ of
 2973 the targets using the Turner1999 energy parameter sets. Using
 2974 a similar technique, RNAinverse[113], one of the oldest inverse
 2975 folding tools included in the ViennaRNA package, uses an adap-
 2976 tive random walk to minimize base-pair distance. The distance
 2977 is computed by comparing the MFE structure of the mutated se-
 2978 quence with the target structure. In addition, RNAinverse allows
 2979 for designing more probable sequences using the partition func-
 2980 tion optimization. The latter allows for more stable designed
 2981 sequences that mostly fold into MFE structures different from the
 2982 target structure. On an attempt to improve RNAinverse, many
 2983 other tools have been suggested INFO-RNA [18], RNA-SSD [5] and
 2984 DSS-Opt [126]. The most recent tools also include RNAPOND [217]
 2985 and MaiRNAlFold [130].

2986 antaRNA [100] is also a recent program available since 2015,
 2987 and it provides a web server for friendly usability. It utilizes an
 2988 *ant-colony* optimization, in which an initial sequence is generated
 2989 via a weighted random search, and the *fitness* of that sequence is
 2990 then used to refine the weights and improve subsequences over
 2991 generations. It provides many other interesting features, such
 2992 as the sequence and target GC-content constraints. It also pro-
 2993 vides a fast python script that includes the options from the web
 2994 server presented through a command line. Other tools also pro-
 2995 vide this dual advantage but implement different optimization
 2996 techniques. NUPACK:design [219] uses a tree decomposition tech-
 2997 nique and the ensemble defect as objective function to design
 2998 qualitatively good sequences. incaRNAbinv [44] is a program
 2999 for fragment-based RNA design. incaRNAbinv's web server com-
 3000 bines two complementary methodologies: IncaRNAtion [146] and
 3001 RNAbinv [209]. IncaRNAtion generates a GC-weighted partition
 3002 function for the target structure, and then adaptively samples se-
 3003 quences from it to match the desired GC-content. RNAiFold [60]
 3004 employs constraint programming that exhaustively searches over
 3005 all possible sequences compatible with a given target. RNAiFold
 3006 [60] has the particularity of designing synthetic functional RNA
 3007 molecules.

3008 So far, except for Modena and antaRNA, most of the computation
 3009 tools presented in previous sections do not account for pseudo-
 3010 knotted RNA target structures, which represents a disadvantage,
 3011 knowing their implications in realizing ncRNA biological func-
 3012 tions. The following section reviews existing RNA design tools
 3013 that support pseudoknotted secondary structures.

3014 4.4.2 Pseudoknotted RNA inverse folding tools

3015 Designing RNA sequences for pseudoknotted targets is computa-
 3016 tionally more expensive than pseudoknot-free targets. For that
 3017 reason, many of the studies addressing the inverse folding of RNA
 3018 considered only pseudoknot-free secondary structures. There are,
 3019 however, some exceptions: MCTS - RNA [215], antaRNA[100], Modena
 3020 and Inv[59]. The computation tool presented in Chapter 5 of our
 3021 work also considers pseudoknots. This section gives an overview
 3022 of each of these tools.

3023 Inv was one of the first inverse folding tools handling pseudo-
 3024 knotted RNA target structures, but it was restricted to a specific
 3025 type of pseudoknot pattern called 3-crossing nonplanar pseudo-
 3026 knots.

3027 More recently, MCTS-RNA's authors suggested a new technique
 3028 that deals with a broader type of pseudoknots. It uses a monte
 3029 carlo tree search ([MCTS](#)) technique which has recently shown
 3030 exceptional performance in Computer Go. The [MCTS](#) allows ini-
 3031 tialising a set of [RNA](#) sequence solutions in MCTS-RNA and the
 3032 solutions are further improved through local updates at the nu-
 3033 cleotide positions.

3034 Another approaches (Modena, antaRNA) implements different
 3035 strategies one which is a multi-objective ant-colony optimisa-
 3036 tion and the another one which is a multi-objective evolutionary
 3037 algorithm. Although the first versions were implemented for
 3038 pseudoknot-free structure [100, 188], they have since been ex-
 3039 tended to support pseudoknotted [RNAs](#) [99, 189].

3040 Each of the tools mentioned above rely on a folding tools
 3041 that predicts pseudoknotted secondary structure: MCTS-RNA uses
 3042 pkiss whereas the other tools (antaRNA and Modena) support two
 3043 folding tools such as HotKnots and IPKnot. In the context of this
 3044 work, two folding tools are used HotKnots and IPKnot, and they
 3045 support the two main types of pseudoknot patterns (i.e. H-type
 3046 and K-type) contained in the benchmark data used to evaluate
 3047 our result in [Chapter 5](#). Both pseudoknotted and pseudoknot-free
 3048 benchmark data sets are considered in this work. The following
 3049 section describes the benchmark data used to evaluate our pro-
 3050 posed [EA](#) tool.

3051 4.5 BENCHMARKING THE INVERSE FOLDING TOOLS

3052 The validation of the designed [RNA](#) sequences using computa-
 3053 tional methods often requires biological experiments. Because of
 3054 the high cost of experimental techniques, most investigators limit
 3055 their guarantee to using benchmark datasets [23] in general. For
 3056 pseudoknot-free design tools, two benchmark datasets are mostly
 3057 used in the literature—(i) RFAM¹: a collection of [RNA](#) families,
 3058 each represented by multiple sequence alignments, consensus
 3059 secondary structures and covariance models—(ii) Eterna100 [3]:
 3060 a collection of hundred [RNA](#) secondary structures extracted from
 3061 the EteRNA Puzzle game². For [RNA](#) inverse tools that support
 3062 pseudoknots, the PseudoBase++[191] dataset is often considered.
 3063 This section provides references, descriptions and the cleanup
 3064 procedure applied for the three data sets mentioned above.

¹ The Rfam database <https://rfam.xfam.org/>

² The EteRNA game <https://eternagame.org/>

3065 The Eterna100 dataset [106] is available in two versions and
3066 both contain a set of 100 target structures extracted from the
3067 EteRNA puzzle game and classified by their degree of difficulty.
3068 The Eterna100-V1 was initially designed using ViennaRNA 1.8.5,
3069 which relies on Turner1999 energy parameters [197]. Out of
3070 the 100 target secondary structures, 19 turned out to be unsolv-
3071 able using the version of ViennaRNA 2.4.14 (which relays on the
3072 Turner2004 [123]). Subsequently, an Eterna100-V2 [106] was
3073 released in which the 19 targets were slightly modified to be solv-
3074 able using ViennaRNA 2.4.14 and any version that supports the
3075 Turner2004 energy parameters. The main difference between the
3076 two dataset relay on the energy parameters used to generate the
3077 data.

3078 The non-EteRNA (a subset of the RFAM) dataset in a set of 63
3079 experimentally synthesized targets that Garcia-Martin et al. [60]
3080 recently used to benchmark a set of ten inverse folding algorithms,
3081 which from our knowledge, is the most recent and comprehensive
3082 benchmark of current state-of-the-art methods. The dataset is
3083 collected from 3 sources: the first dataset called **dataset A** which
3084 contains 29 targets collected from RFAM and also used in [47,
3085 188] and the second called **dataset B** is a collection of 24 targets
3086 used in [47] and added to that the 10 structures used in [171].

3087 The PseudoBase++ is a set of 266 pseudoknotted RNA structures
3088 used to benchmark Modena. It was initially 342 RNA secondary
3089 structures, but because of the redundancy and the non-canonical
3090 base pairs 76 structures were excluded. To group the dataset with
3091 respect to the pseudoknot motifs, we used the test data from
3092 antaRNA's paper. The test data contains 249 grouped into four
3093 categories: 209 hairpin pseudoknots (H), 29 bulge pseudoknots
3094 (B), 8 complex hairpin pseudoknots (cH) and 3 kissing hairpin
3095 pseudoknots (K). Out of the 266 structures, only 185 (with 150 H-
3096 type, 3 K-type, 25 B-type and 7 cH-type) structures were included
3097 in the test data. So for that reason, we have used only 185 target
3098 structures for the pseudoknot motif performance comparison and
3099 the 266 structures for the different target lengths performance
3100 comparison.

3101 When the benchmark datasets rely on a particular energy pa-
3102 rameter set, the performance of a given inverse RNA folding tool
3103 evaluated on these datasets will also be related to the choice of
3104 the RNA folding tool's energy parameter set. If the benchmark
3105 datasets do not rely on a particular energy parameter set, the ro-
3106 bustness of the inverse RNA tool will be its capability to perform
3107 well on different energy parameter sets.

3108 4.6 CONCLUSION

3109 In summary, the RNA inverse folding problem is still computa-
3110 tionally challenging because there are many objective functions
3111 and different ways of evaluating computational tools. Solving this
3112 problem is particularly interesting in RNA synthetics, RNA nanos-
3113 tructure design, and emerging fields such as bioengineering. We
3114 presented a comprehensive literature review of existing compu-
3115 tational methods that addressed this problem in this chapter. The
3116 existing approaches have some advantages and disadvantages,
3117 depending on the techniques implemented. NUPACK for example—
3118 despite its well-defined objective function—still has difficulty
3119 designing sequences for large targets and most of the EteRNA100
3120 targets. In contrast, ERD because of its powerful decomposition
3121 method, which allows dealing quickly with large targets (On
3122 RFAM 1.0 with target's length between 400 – 1400) but is still a
3123 big challenge to solve more than 65% of the EteRNA100-V2 using
3124 the Turner2004 energy parameter sets. On another side, NEMO,
3125 one of the most recent tools, can solve more than 90% of the
3126 EteRNA100-V1 dataset using an old version of ViennaRNA pack-
3127 age, which is based on Turner1999 energy parameter sets [197].
3128 The sentRNA's machine learning model also relied on the same old
3129 version of ViennaRNA package and, by adding a refinement on the
3130 machine learning model, sentRNA solves 78% of EteRNA100. With-
3131 out this refinement, sentRNA can only solve 48% of EteRNA100's
3132 targets, which can represents another limitation. For the EAs ERD
3133 and MODENA, none of them can solve more than 65% of EteRNA100
3134 using the Turner2004 energy parameter sets.

3135 In the next chapter, we will introduce a simple evolutionary
3136 algorithm called aRNAque that implements a Lévy mutation and
3137 allows significant improvements to the existing tools.

3138

3139 AN EVOLUTIONARY ALGORITHM FOR INVERSE 3140 FOLDING INSPIRED BY LÉVY FLIGHTS.

3141 In the previous chapter of our work, we presented the **RNA** de-
3142 sign as an optimization problem and provided a significant lit-
3143 erature review on the existing tools addressing that problem.
3144 We highlighted some limitations of the existing tools, particu-
3145 larly those implementing evolutionary algorithms. One of the
3146 main challenges of evolutionary algorithms is to avoid decep-
3147 tion, which is the fast convergence to a local optimum. Most **EAs**'
3148 early convergence to a local optimum is due to the local search
3149 implementation, which is the consequence of the local mutation
3150 scheme.

3151 To avoid this pitfall, an alternative mutation scheme to the
3152 classical local search is the Lévy mutation. We propose an evo-
3153 lutionary algorithm that implements a similar Lévy mutation in
3154 this chapter but adjusts to the **RNA** design problem. This mutation
3155 scheme is focused on local search but also searches at all other
3156 scales to avoid becoming trapped. its long-range search permits
3157 designing **RNA** sequences of higher positional entropy. Our im-
3158 plementation, called **aRNAque** is available on GitHub as a python
3159 script. Compared to existing inverse folding tools, the benchmark
3160 results show improved performance on both pseudoknot-free
3161 and pseudoknotted datasets. Much of materials in this chapter
3162 has been previously published in [128, 129].

3163 5.1 MATERIAL AND METHODS

3164 This section provides a detailed description of **aRNAque** algorithm
3165 in general and in particular the Lévy mutation scheme imple-
3166 mented.

3167 5.1.1 *aRNAque's mutation operator*

3168 The previous chapters provided an overview of **EA**, emphasizing
3169 its application to the **RNA** inverse folding problem. One of
3170 the essential components of **EAs** is the mutation operator. Our
3171 tool, **aRNAque**, implements a simple **EA** that uses a Lévy muta-

tion to explore at different scales the solution space. In addition, our mutation allows explicitly controlling the GC-content of the designed RNA sequences. This section presents in detail our proposed mutation operator.

For a given target RNA secondary structure in its string representation σ^* of length L , the space of potential solutions to the inverse folding problem is $\{A, C, G, U\}^L$. An evolutionary algorithm explores the space of solutions through its move (or mutation) operator. To explore the search space of compatible sequences (sequences with canonical base pairs at the corresponding open and closed bracket positions) with σ^* exclusively, we propose a mutation step that depends on the nucleotide canonical base pair probability distribution.

Let $\mathcal{N} = \{A, C, G, U\}$ be the set of nucleotides weighted respectively by the probabilities:

$$P_{\mathcal{N}} = \{w_A, w_C, w_U, w_G\}$$

and $\mathcal{C} = \{AU, UA, CG, GC, UG, GU\}$ be the set of canonical base pairs weighted respectively by the probabilities:

$$P_{\mathcal{C}} = \{w_{AU}, w_{UA}, w_{CG}, w_{GC}, w_{UG}, w_{GU}\}$$

where

$$\sum P_{\mathcal{N}} = 1, \sum P_{\mathcal{C}} = 1$$

Our evolutionary algorithm relies on the flexibility of the mutation parameters $P_{\mathcal{N}}, P_{\mathcal{C}}$. These parameters allow explicit control of the GC-content of the RNA sequences during the designing procedure.

We examined the binomial and Zipf distributions:

- Binomial mutation: here U has a binomial distribution:

$$P(U = n) = \binom{L}{n} \mu^n (1 - \mu)^{L-n}$$

for some $0 \leq \mu \leq 1$, such that $u = \mu \cdot L$. We can think of this mutation mode arising from each nucleotide of an RNA sequence independently undergoing a point mutation with probability μ , i.e. μ is the per-nucleotide or point mutation rate.

- Lévy mutation: U has a Zipf distribution given by:

$$P(U = n) = \frac{1/n^c}{\sum_{k=1}^L 1/k^c}$$

3195 where $c > 0$ is the value of the exponent characterizing the
 3196 distribution.

3197 **Figure 5.1** shows the distribution of the number of point muta-
 3198 tions on a sequence of length 88 nucleotides for both mutation
 3199 schemes. Both distributions have the same mean, and the differ-
 3200 ence between the two distributions is more perceptible on their
 3201 tails.

3202 In the rest of this work, a local mutation will refer to a binomial
 3203 mutation with parameter $\mu \approx 1/L$.

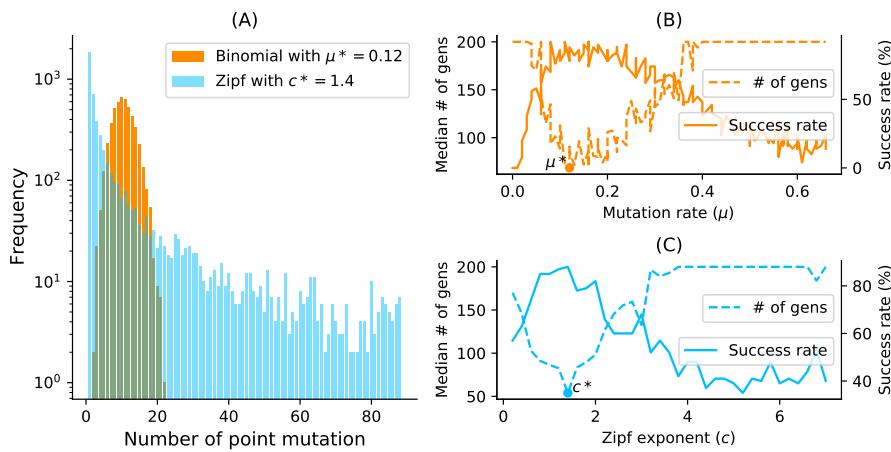


Figure 5.1: **Binomial vs. Zipf distributions.** (A) Samplings Binomial and Zipf distributions for the best binomial mutation rate μ^* (respectively c^* for the best Zipf exponent parameter). Both distributions have a mean of 8.7 point mutations for a sequence of length 88 nucleotides. (B) Tuning of binomial mutation rate parameter. For each $\mu \in [0, 1]$ with a step size of 0.005 and the pseudoknotted target PKB00342 of length 88, 50 sequences were designed using aRNAque. (B) shows the median generations and the success percentage *vs.* the mutation rate (μ). The best mutation rate is $\mu^* = 0.085$ (with a median number of generation 93.5 and a success rate of 92%). (C) Tuning of Levy exponent. Similar to (B), for each $c \in [0, 7]$ with a step size of 0.1 and for the same pseudoknotted target structure, 100 sequences were designed using aRNAque. It shows the median generations and the percentage of success *vs.* the exponent parameter (c). The Zipf exponent distribution that produced the highest success rate and the minimum number of generations is $c^* = 1.4$.

3204 We present the mutation algorithm in [algorithm 1](#). This muta-
 3205 tion algorithm is intergraded in a unified EA framework, allowing

3206 to update RNA sequence solution at each iteration or generation.
 3207 After we apply the mutation operation to the population of RNA
 3208 sequences, we evaluate the newly generated population; this is
 3209 usually done using an objective or fitness function. In the follow-
 3210 ing section, we describe the different objective functions taken
 3211 into account in our implemented EA.

3212 **5.1.2 aRNAque's objection functions**

3213 Our EA reaches its performance through the minimization of
 3214 three objective functions:

- Hamming distance from the target structure: Since the main goal of the inverse folding problem is to find sequences that fold into a given target secondary structure σ^* , the simple fitness measurement f of an RNA sequence ϕ can be defined as follows:

$$f(\phi, \sigma^*) = \frac{1}{1 + d_h(\sigma^{MFE}(\phi), \sigma^*)} \quad (5.1)$$

3215 where $d_h(\cdot, \cdot)$ is the hamming distance on the structure
 3216 space (structures are in dot and bracket representation)
 3217 defined in [Equation 1.12](#).

- normalized energy distance ([NED](#)): It is used to minimize the free energy of the designed sequences. (See [Equation 1.14](#))
- ensemble defect ([ED](#)) [[219](#)]: Here, we use the [ED](#) as a second objective function for refinement after having at least one sequence that folds into the target in the current population. It is defined in [Equation 1.13](#).

3225 To minimize the [NED](#) and the hamming distance of a population of
 3226 RNA sequences, instead of combining both measurements to form
 3227 a multi-objective function, we use them separately at a different
 3228 level of our EA. We use the [NED](#) as a selection weight for the
 3229 sequences that will be mutated, and the hamming distance is
 3230 used as a weight to elite ten best sequences that will always move
 3231 to the next generation. Therefore the selection method we use
 3232 is the *fitness proportionate selection*, also known as roulette wheel
 3233 selection [[110](#)]. Once we have at least one sequence that folds into
 3234 the given target in the current population (for the successful case),
 3235 we start a random walk in its neutral network by minimizing

Algorithm 1: aRNAque's mutation algorithm

```

/*  $P' = \{\phi'_1 \dots \phi'_n\}$ : the mutated population;
 $P = \{\phi_1 \dots \phi_n\}$ : a list of  $n$  RNA sequences to mutate;
 $P_C = \{w_{AU}, w_{UA}, w_{CG}, w_{GC}, w_{UG}, w_{GU}\}$ : a vector containing the
weights associated with each canonical base pairs;
 $P_N = \{w_A, w_U, w_C, w_G\}$ : a vector containing the weights
associated with each nucleotide;
 $\mathcal{D}$ : a given probability distribution (Lévy or Binomial) with
parameter  $p$  and  $L$ . Where  $L$  is the length of the target
RNA structure */
```

Input: $P, \mathcal{D}(p, L), P_C, P_N$

Output: P'

- 1 $\{B_i\} \sim \mathcal{D}(p, L)$, where $i \in \{1, 2, \dots, n\}$; // Draw n random numbers that follows a given distribution $\mathcal{D}(p, L)$ (Lévy or Binomial). B_i is the number base pairs to mutate
- 2 $\{U_i\} \sim \mathcal{D}(p, L)$, where $i \in \{1, 2, \dots, n\}$; // Draw n random numbers that follows the same distribution as B_i (Lévy or Binomial). U_i is the number non base pair positions to mutate
- 3 **for** $i \in \{1, 2, \dots, n\}$ **do**
 - 4 $\phi' \leftarrow P_i$; // Assign the sequence $\phi_i \in P$ to ϕ'
 - 5 **for** $j \in \{1, 2, \dots, U_i\}$ **do**
 - 6 $r \in \{1, 2, \dots, L\} \sim \mathcal{U}$; // select uniformly a random position in the RNA sequence ϕ'
 - 7 $n_j \in \{A, U, C, G\} \sim P_N$; // select a random nucleotide n_j with respect to P_N
 - 8 $\phi'_r \leftarrow n_j$; // replace the nucleotide at position j in the RNA sequence ϕ' with n_j
 - 9 **for** $j \in \{1, 2, \dots, B_i\}$ **do**
 - 10 $k_j \in \{AU, UA, CG, GC, UG, GU\} \sim P_C$; // select a random base pair k_j with respect to P_C
 - 11 $b \in \{(b_1, b_2)_i\} \sim \mathcal{U}$; // select uniformly a random pair of base pair positions
 - 12 $\phi'_b \leftarrow k_j$; // replace respectively the nucleotides at the base pair position $b_i \in b$ by $k \in k_j$
 - 13 $P' \leftarrow P' \cup \phi'$; // Add ϕ' to the list P'

3236 the ensemble defect function (Equation 1.13). The next section
 3237 provides more detailed information about the core of our EA and
 3238 the full pseudo-code.

3239 Now that we have defined the mutation and selection operators
 3240 implemented in our EA, we will describe the general algorithm
 3241 in the following section.

3242 **5.1.3 aRNAque's EA**

3243 As described in the introductory chapter, an EA starts with an ini-
 3244 tial population of solutions and sequentially applies the mutation
 3245 and selection operators on the solutions through generation until
 3246 a termination criterion is satisfied. How does aRNAque generate
 3247 the initial population of RNA sequences? This section describes
 3248 how the initial population is generated and then provides the
 3249 core pseudocode of our EA.

3250 For a given population size n and a target structure \mathcal{S}^* of length
 3251 L , an initial population P is generated randomly as follows:

- 3252 1. Select randomly L nucleotides in \mathcal{N}
- 3253 2. Identify the base pair position (i, j) in the random sequence,
 3254 select randomly a base pair in the set of canonical base pairs
 3255 \mathcal{C} and fix the first nucleotide of the selected canonical base
 3256 pair at the position i and the second at position j .
- 3257 3. Repeat 2. for all base pair positions in the target structure
- 3258 4. Repeat 1. 2. and 3. n -times.

3259 Let T be the maximum number of generations and F_t the set of
 3260 all sequences found at a given time t . After the initial population
 3261 of RNA sequences is generated, our algorithm is described in
 3262 **algorithm 2**. The stopping criteria are two: 1) the number of
 3263 generations (t) is equal to the max number of generations (T)
 3264 or 2) the minimum hamming (or base pair) distance of the best
 3265 RNA sequence solution to the target is 0 (i.e the maximum fitness
 3266 value is 1).

3267 In sum, our EA relies on three objective functions and imple-
 3268 ments a Lévy flight mutation scheme. We assess the performance
 3269 of our EA and the existing tools using three benchmark data
 3270 sets presented in the previous chapter (Section 4.5). The follow-
 3271 ing section describes the benchmark protocols applied for each
 3272 data set and different RNA inverse folding tools considered in the

Algorithm 2: aRNAque' EA

```

/*  $P' = \{\phi'_1 \dots \phi'_n\}$ : the best population;
 $P = \{\phi_1 \dots \phi_n\}$ : the initial population of  $n$  RNA sequences;
 $P_C = \{w_{AU}, w_{GU}, w_{GC}\}$ : a vector containing the weights
associated with each base pair;
 $P_N = \{w_A, w_U, w_C, w_G\}$ : a vector containing the weights
associated with each nucleotide;
 $\mathcal{D}$ : a given probability distribution (Lévy or Binomial) with
parameter  $p$  and  $L$ , where  $L$  is the length of the target RNA
structure;
 $T$ : the maximum number of generations;
 $n$ : the population size ;
 $f(\cdot)$ : the fitness function used. It can be the hamming,
base-pair or energy distance;
 $\sigma^*$ : the target structure in its string representation;
 $\mathcal{P}$ : the energy parameters used for the folding */
```

Input: $n, T, P_N, P_C, P, \mathcal{D}(p, L), f(\cdot), \sigma^*, \mathcal{P}$

Output: Best population P'

```

1  $P' \leftarrow P$ ; // Assign the initial population to the best
population
2  $t \leftarrow 0$ ; // Initialize the number of generations to 0
3 while  $t \leq T$  &  $f(\sigma^{MFE}(\phi_b), \sigma^*) \neq 1$  do
4    $\Sigma \leftarrow \{\arg \min_{\sigma \in \Sigma} \Delta G(\phi_i, \sigma, \mathcal{P})\};$  // Fold each sequence
     $\phi_i \in P'$  and store them in  $\Sigma$ . Where  $i \in \{1, 2, \dots, n\}$ ,  $\Gamma$ 
    the structural ensemble and  $\Delta G(\phi_i, \sigma)$  the free energy
    computed using the parameters  $\mathcal{P}$ 
5    $\kappa = \lfloor (n \times 0.1) \rfloor$ ; // The number of RNA sequences to copy
    in the next generation without mutating them.
6    $F \leftarrow \{f(\sigma, \sigma^*) | \forall \sigma \in \Sigma\}$ ; // Evaluate the fitnesses of
    the folded population to the target strucre  $\sigma^*$  and
    store them in a list  $F$ 
7    $E_\kappa \leftarrow \{\phi_1 \dots \phi_\kappa\} \sim F$ ; // copy of the 10% best sequence
    based on their fitnesses  $F$ .
8    $P_S \leftarrow \{\phi_i\} \sim F$ , where  $i \in \{1, 2, \dots, n - \kappa\}$ ; // Randomly
    sample  $(n - \kappa)$  RNA sequences from  $P'$  with respect to
    their fitnesses  $F$ .
9    $M \leftarrow \text{mutate}(P_S, \mathcal{D}(p, L), P_C, P_N)$ ; // Mutated the
    selected sequences using the mutation algorithm
    presented in the main text in out paper.
10   $P_b \leftarrow M \cup E_\kappa$ ; // Combine the mutated population and the
    best solutions to form the new population that will be
    evolved in the next generation
11   $\phi_b \leftarrow \arg \max_{\sigma \in \Sigma} f(\sigma, \sigma^*)$ ;
12   $t \leftarrow t + 1$ ; // Increment the time step (the number of
generations)

```

3273 context of this work. Furthermore, it provides an overview of
 3274 various folding tools and the configuration parameters used for
 3275 the benchmark.

3276 *5.1.4 Benchmark parameters and protocols*

3277 For the benchmark results presented in this work, we use three
 3278 datasets: the Eterna100 dataset, RFAM dataset and PseudoBase++
 3279 dataset. Depending on the datasets, a specific RNA folding tool
 3280 is used. This section gives more details about aRNAque's parame-
 3281 ters, energy parameters and other tools parameters used for the
 3282 benchmark results presented in this chapter.

3283 *Folding tools*

3284 Two tools for pseudoknotted RNA folding are considered in this
 3285 work: HotKnots and IPknot. For pseudoknot-free RNA folding, we
 3286 used RNAfold. For the mutation parameter and GC-content anal-
 3287 ysis presented in our work, we used IPknot, and both HotKnots
 3288 and IPknot for PseudobBse++ benchmarks. To be able to use
 3289 HotKnots in aRNAque without copying aRNAque in the bin direc-
 3290 tory of Hotknots, we have performed some modifications on
 3291 Hotknots source code. Details on the modifications are provided
 3292 in the [Section B.5](#). Furthermore, we considered pkiss, a well
 3293 known tool for K-type pseudoknot prediction, but since the PseudoBase++
 3294 dataset contains just 4 K-type pseudoknotted structures and pKiss
 3295 has higher time complexity ($O(n^6)$), we did not find it efficient
 3296 for the benchmark we performed.

3297 *Mutation parameters tuning*

3298 The main challenge for an evolutionary algorithm is to find opti-
 3299 mum parameters such as mutation rate, population size and se-
 3300 lection function. We used 80 pseudoknotted targets with lengths
 3301 from 25 to 181 nucleotides for the mutation parameter analysis.
 3302 We set the maximum number of generations T to 200 and the
 3303 population size n to 100. The stopping criteria are two: 1) the
 3304 number of generations (t) is equal to the max number of genera-
 3305 tions (T) or 2) the minimum hamming (or base pair) distance of
 3306 the best RNA sequence solution to the target is 0. The best muta-
 3307 tion parameters (c^* for Levy and μ^* for Binomial) are those that
 3308 have the lowest median number of generations. The best muta-
 3309 tion parameters obtained for both binomial and Lévy mutation

3310 modes are used to benchmark and compare the results on the
 3311 entire datasets of RNA structures.

3312 *Benchmark on the PseudoBase++ dataset*

3313 Four benchmarks are performed on the pseudoknotted dataset:
 3314 1) mutation parameter analysis, 2) the GC-content and diversity
 3315 analysis, 3) Local search versus Lévy search, 4) aRNAque (Lévy
 3316 search) versus antaRNA. For the aRNAque (Binomial and Lévy)
 3317 case, the four benchmarks share the same number maximum
 3318 number of generations ($T = 200$), population size ($n = 100$),
 3319 stopping criteria ($t = T$ or min fitness equals 0). For the antaRNA
 3320 benchmark, the maximum number of iterations was set to 1200,
 3321 and a slight modification was made to allow the support of the
 3322 folding tool HotKnots (See [Section B.5](#)). For booth tools and each
 3323 benchmark, 20 runs were launched independently in parallel on
 3324 a computer with the same resources, resulting in 20 designed
 3325 sequences per pseudoknotted target structure. To measure the
 3326 performance of each tool, each designed sequence s is folded into
 3327 a secondary structure \mathcal{S} and the similarities between \mathcal{S} and \mathcal{S}^* are
 3328 computed using the base pair distance. For the GC-content bench-
 3329 mark, four GC-content values are considered, $\{0.25, 0.5, 0.75, 1\}$
 3330 and the setting of each tool remains the same.

3331 *Benchmark on the Eterna100 dataset*

3332 We performed two benchmarks are one the Eterna100 dataset: 1)
 3333 a benchmark on the Eterna100-V1 dataset using the Turner1999
 3334 energy parameter and the both versions of aRNAque (one point
 3335 and Lévy mutation), 2)a benchmark on the Eterna100-V2 dataset
 3336 using the Turner2004 energy parameter and both versions of
 3337 aRNAque (one point and Lévy mutation). For each of the Eterna100
 3338 benchmark we used the same evolutionary algorithm parameters;
 3339 a maximum of $T = 5000$ generations (i.e. a maximum of 500,000
 3340 evaluations), a population size of $n = 100$ and the same stopping
 3341 criteria (the number of generation $t = T$ or min fitness equals
 3342 0). For both local and Lévy search, 5 runs were launched inde-
 3343 pendently, which results in 5 designed sequences per target. We
 3344 define success rate simply as the number of successfully designed
 3345 targets. A target is considered successfully designed when at least
 3346 one of the designed sequences folds into the target structure.

3347 For the benchmarks peformed on ERD, NUPACK, and SentRNA the
 3348 default parameters were used. For NEMO, the number of iteration

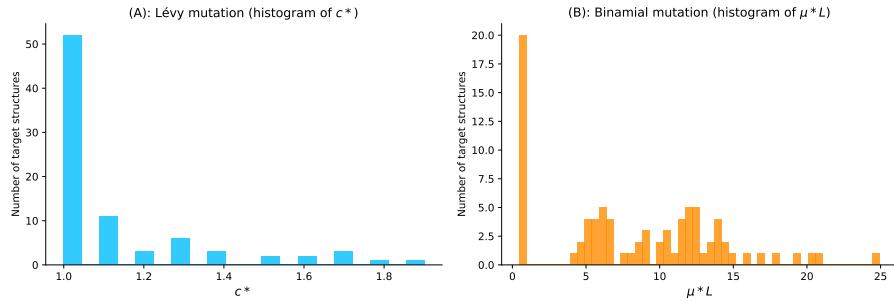


Figure 5.2: Parameter tuning for both binomial and Lévy mutation schemes. (A) Lévy mutation parameter tuning. Histogram of best exponent parameter (c^*) for a set of 81 target structures with different pseudoknot patterns and various lengths. The most frequent best exponent value is 1. (B) Binomial parameter tuning. Histogram of best mutation rate (μ^*) for the same set of 81 target structures with different pseudoknots and various lengths. The most frequent best parameter is the low mutation rate ($\approx 1/L$). For some structures, the best mutation rate is the high one for different lengths as well.

3349 was set to 2500 and for RNAinverse the objective function was set
 3350 to be the partition function and the number of iteration at 1200.

3351 *Benchmark on the non-Eterna100 dataset*

3352 For the non-EteRNA dataset, only the Turner2004 energy param-
 3353 eters were used. The maximum number of generations was set
 3354 to be 5000. The mutation parameters (P_C and P_N) were chosen
 3355 to be close to the nucleotide distribution of the RNA sequence in
 3356 nature [47].

3357 **5.2 EXPERIMENTAL RESULTS**

3358 As mentioned in Chapter 4, the validation of computational tools
 3359 for RNA inverse folding can include in vivo or in vitro experiments.
 3360 In the context of this work, only in silico experiments are used to
 3361 evaluate the performance of the existing tools, including aRNAque.
 3362 This is done through a benchmark protocol described in the pre-
 3363 vious section. This work’s computational tools require an RNA
 3364 secondary structure as an input target. Two RNA secondary struc-
 3365 tures are considered: the pseudoknot-free and the pseudoknotted
 3366 target structures, and both are supported in aRNAque. Therefore,
 3367 we evaluate aRNAque performance for different target secondary

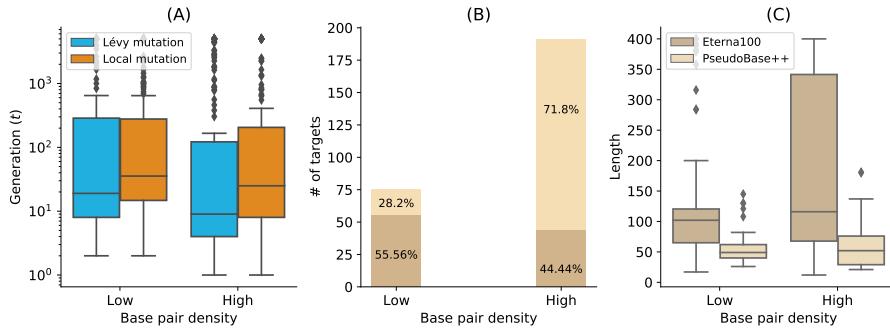


Figure 5.3: Lévy mutation vs. Local mutation: performance analysis with respect to the base-pair density. The higher the base-pair density is, the more useful the Lévy mutation scheme to speed up the optimization EA. (A) Distributions of number of generations for the low and high base-pair density targets of the Eterna100 dataset. (B) Percentages of targets with low and high base-pair density for the Eterna100 and PseudBase++. (C) The length distributions of the low and high base-pair density pseudoknot-free and pseudoknotted targets.

3368 structures separately. Three data sets of secondary structure tar-
 3369 gets are used: the Eterna100 and non-Eterna100 that contain
 3370 pseudoknot-free targets, and the PseudoBase++ which contains
 3371 only pseudoknotted targets. Using the three data sets, this section
 3372 presents the experimental results concerning the quality of the
 3373 designed RNA sequences (i.e. the GC-content and diversity or
 3374 positional entropy), the CPU required and the success rate of each
 3375 tool considered when benchmarked using a specific data set.

3376 5.2.1 aRNAque's performance on pseudoknot-free target structures

3377 We compared the performance of aRNAque for pseudoknot-free
 3378 target using the benchmark datasets: the non-Eterna100 and
 3379 the Eterna100. This subsection presents the statistical results
 3380 obtained compared to benchmarked existing tools and the results
 3381 found in the literature. In addition, we compared the performance
 3382 of aRNAque (Lévy mutation) to the one of Ivry et al. [88] on a
 3383 tripod pseudoknot-free RNA secondary structure.

3384 5.2.1.1 Performance on Eterna100 dataset

3385 A first benchmark was performed on the Eterna100 datasets. First,
 3386 on the Eterna100-V1 dataset, the Lévy flight version of aRNAque

Table 5.1: Summary of performance of aRNAque vs the 7 other algorithms benchmarked on EteRNA100-V1 by Anderson-Lee et al. [3] (using the recent energy parameter sets, the Turner2004)

Methods	Number of puzzles solved
aRNAque	72/100
RNAinverse	66/100
Learna	66/100
ERD	65/100
SentRNA, NN + full moveset	60/100
MODENA	54/100
NEMO	50/100
INFO-RNA	50/100
NUPACK	48/100
DSS-Opt	47/100
RNA-SSD	27/100

3387 successfully designed 89% of the targets and the one-point mutation
 3388 (local mutation) version achieved 91% of success, suggesting
 3389 that for some target structures, local mutation can outperform the
 3390 Lévy mutation scheme. Combining the two solutions, aRNAque
 3391 solved in total 92% of the targets of Eterna100-V1.

3392 When analysing the performance of Lévy flight for low and
 3393 high base pair densities separately, the median number of genera-
 3394 tions of high base pair density targets was lower than the one
 3395 with low base-pair density (8 generations for high density and
 3396 18 for the low base pairs density targets). The same observation
 3397 was drawn for the success rate. For the low base-pair density
 3398 targets, the Lévy flight achieved 87% (49/56) success whereas,
 3399 for the high base-pair density, it achieved 91% (40/44). The same
 3400 analysis can be done when comparing the one-point mutation
 3401 results for the high-density targets to the Lévy flight mutation.
 3402 The median number of generations for the low-density targets
 3403 when using a one-point mutation operator was 34 (respectively
 3404 24 for the high base pair density targets) (see Figure 5.3A).

3405

3406 Another benchmark was performed on Eterna100-V2 with aRNAque
 3407 achieving a 93% success rate when combining the designed solu-
 3408 tions for both mutation schemes. Compared to recently reported

Table 5.2: **Summary of performance of aRNAque vs the 10 other algorithms benchmarked on the non-EteRNA100** by Anderson-Lee et al. [3]

Methods	Number of puzzles solved
SentRNA, NN + full moveset	57/63
ERD	54/63
SentRNA, NN + GC pairing	53/63
SentRNA, NN + All pairing	53/63
aRNAque	52/63
RNA-SSD	47/63
SentRNA, NN only	46/63
INFO-RNA	45/63
MODENA	32/63
NUPACK	29/63
IncaRNAtion	28/63
Frnakenstein	27/63
RNAinverse	20/63
RNAfbinv	0/63

3409 benchmark results [106], aRNAque achieved almost similar per-
 3410 formance to NEMO on Eterna-V2: one target was unsolved by all
 3411 existing tools and one target solved only by NEMO remained un-
 3412 solved by aRNAque, outperforming all existing EA methods.

3413

3414 For the robustness analysis, Table 5.1 presents the benchmark
 3415 results on Eterna100-V1 using the Turner2004 energy parame-
 3416 ters sets. It shows that the evolution algorithm we propose can
 3417 solve $\approx 72\%$ of the dataset, and it surpasses the 4 methods we
 3418 benchmarked and all the tools already benchmarked in [171].
 3419 We can also solve approximately 23 targets more than NUPACK,
 3420 which is also minimizing the ensemble defect and that shows the
 3421 importance of a population-based algorithm. Compared to the
 3422 existing EA-based algorithms, our EA can solve approximately 18
 3423 targets more than MODENA and 7 targets more than ERD.

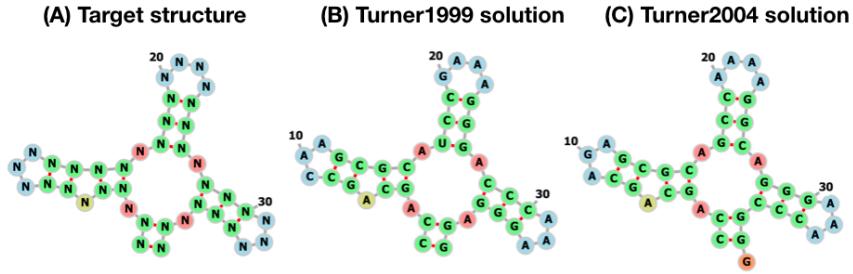


Figure 5.4: **aRNAque’s performance on a TRIPOD secondary structure.**

(A) The tripod target structure. (B) aRNAque’s solution using the Turner1999 energy parameter sets. (C) aRNAque’s solution using the Turner2004 energy parameter sets.

3424 5.2.1.2 *Perfomance on non-Eterna100*

3425 Additionally to the Eterna100 dataset, we also used the non-
 3426 EteRNA dataset collected from the RFAM database to assess the
 3427 aRNAque’s performance on pseudoknot-free target secondary struc-
 3428 ture. Compared to other tools, the statistical results are presented
 3429 in [Table 5.2](#).

3430 The results show that our method surpasses 8/10 of other
 3431 tools. ERD solved 2 more targets than our method because of its
 3432 strong decomposition capacity, which allows it to solve the entire
 3433 **dataset B**. With the advantage that our evolutionary algorithm
 3434 also allows us to fit the nucleotide distribution parameters taken
 3435 from natural **RNA** directly in the mutation parameters, we can
 3436 solve 21/24 targets from the **dataset B**. For the **dataset A** aRNAque
 3437 solves 24/29 targets which means 2 more than the existing tools
 3438 and for the 10 last targets, it solves 7 targets. Adding all these
 3439 solved targets together, we obtain a result of 52/63 as presented
 3440 in [Table 5.2](#).

3441 5.2.1.3 *aRNAque performance on a tripod secondary structure*

3442 Finally, we performed a benchmark on a tripod target secondary
 3443 structure. The tripod secondary structure was used as a third test
 3444 case in the work of Ivry et al. [88], and it does not contain any
 3445 pseudoknot interactions. It comprises four stems, three of which
 3446 with terminal hairpins, surrounding a multibranch loop (See [Fig-](#)
 3447 [ure 5.4A](#)). The tripod target structure was proved to be very chal-
 3448 lenging, especially because of its multiloop component, which
 3449 is also found in some of the unsolved Eterna100 target struc-
 3450 tures. We perform here, for both energy parameters Turner1999

and Turner2004, 100 independent designs, using a population size of 100 RNA sequences and a maximum of 5000 generations. The mutation parameters used are: $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$, $P_N = \{0.7, 0.1, 0.1, 0.1\}$ and $c = 1.5$. When using the Turner2004 energy parameter set, none of the 100 designed RNA sequences was successful (i.e, 0 sequence folds exactly into the target structure after 5000 generations). Figure 5.4B shows one of the best solutions obtained out of 100 designed sequences when using the Turner2004, the designed sequence folds into a structure at one error base-pair distance from the target structure. In contrast, when using the Turner1999 energy parameters, we successfully designed the tripod secondary structure (See Figure 5.4C). The 100 sequences designed folded exactly into the target structure with an average median number of generations 20. When comparing both solutions to the one obtained in [88], aRNAque (with no need of changing the RNA structure distance) can successfully design the multibranch loop component with one base pair error using the Turner2004 energy parameter whereas RNAinverse (with the DoPloCompare distance) failed to design the multi-branch loop, and the solution was at 2 base-pair distance error.

5.2.2 aRNAque's performance on pseudoknotted target structures

Secondly, we assessed the performance of aRNAque in designing pseudoknotted target secondary structures through intensive benchmark on PseudoBase++ dataset. We then compared the results obtained to the one of antaRNA, using both folding tools Hotknots and IPknot. Furthermore, a comparison between local and Lévy mutations is provided.

5.2.2.1 Best mutation parameter analysis on PseudoBase++: Levy mutation vs. local mutation

The advantage of using a Lévy mutation is its capacity to allow simultaneous search at all scales over the landscape. The search at different scale is often dictated by the exponent parameter of the heavy-tailed distribution. In this first subsection, we analyse for 80 pseudoknotted target structures and for both mutation schemes the distributions of the best mutation parameters.

- Binomial mutation: From Figure 5.1B, the critical range was identified to be from 0 to 0.2 and as μ becomes greater than 0.1, the success rate decreases and the average number

3489 of generations increases. For each of the 80 target struc-
 3490 tures with pseudoknots, 20 sequences were designed for
 3491 $\mu \in [0, 0.2]$ with a step size of $1/L$. [Figure 5.2](#)B shows the
 3492 histogram of the best mutation rate found for each target
 3493 structure. Two main regimes are apparent: one where the
 3494 best mutation rate is very low mutation rate ($\approx 1/L$) and
 3495 another where the high mutation rate is optimal.

- 3496 • Lévy mutation: From [Figure 5.1](#)C, the critical range of c was
 3497 identified to be $[1, 2]$. For $c \in [1, 2]$ and a step size of 0.1,
 3498 an optimum exponent parameter c^* was investigated for all
 3499 the 80 target structures. [Figure 5.2](#)A shows the histogram of
 3500 c^* . Contrary to binomial mutation, the optimum exponent
 3501 parameter does not vary too much ($\forall \mathcal{S}, c^* \approx 1$).

3502 [Figure 5.2](#) shows that when using a Lévy mutation, the optimal
 3503 mutation rate is the same for most target structures. In contrast,
 3504 the optimum binomial mutation rate parameter μ^* mostly varies
 3505 with different targets. Although both mutation schemes (for the
 3506 best mutation parameters) have approximately the same success
 3507 rates, the Lévy flight mutation scheme is more robust to different
 3508 targets.

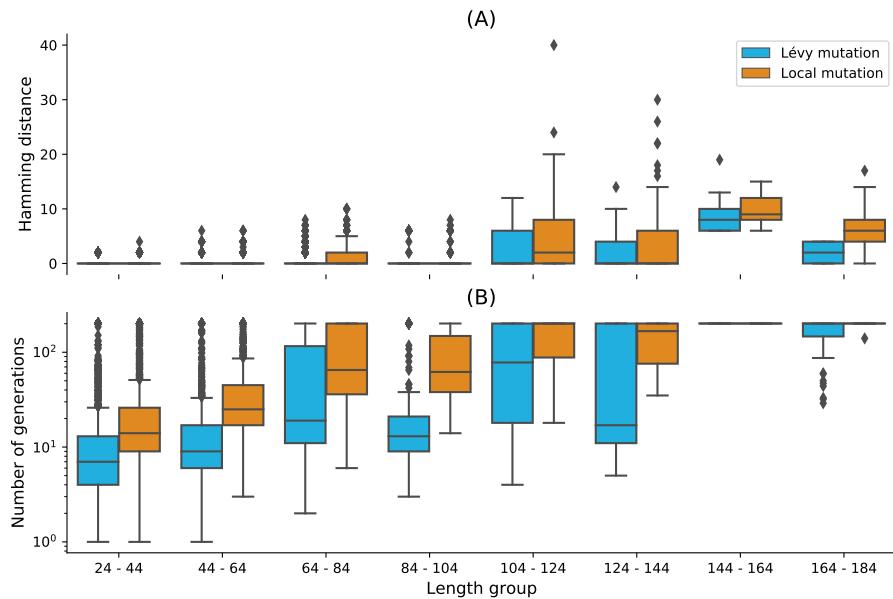


Figure 5.5: Lévy mutation mode vs local mutation (one-point mutation). (A) Hamming distance distributions vs. target structure lengths. (B) Number of generations distributions for different length groups. In both (A) and (B), lower values indicate better performance. The target structures are solvable in less than 100 generations for both mutation schemes and most length groups. Still, the difference in the number of generations gets more significant as the target lengths increase, except for the two last length groups for which both mutation schemes mostly failed. The highest difference in terms of median number of generations is 150 for target lengths in the range [124 – 144] (respectively 123, 49, 46, 16, 7, 0, 0 for the length ranges [84 – 104], [64 – 84], [104 – 124], [44 – 64], [24 – 44], [144 – 164], [164 – 184]). Averaging over all length groups, the median number of generations difference between the Levy mutation and the one point mutation is 48 generations.

3509 5.2.2.2 *Performance on PseudoBase++: Levy mutation vs. local mu-*
 3510 *tation*

3511 **Figure 5.5** shows box plots for the base pair distance (Hamming
 3512 distance) and the number of generations for increasing target
 3513 lengths under our two mutation schemes: binomial at low muta-
 3514 tion rate (or one point mutation) and the Lévy mutation. For each
 3515 pseudoknotted RNA target structure in the PseudoBase++ dataset,
 3516 we designed 20 sequences. The results show that using the Lévy
 3517 mutation instead of a local mutation scheme can significantly in-
 3518 crease the performance of aRNAque. The gain was less significant

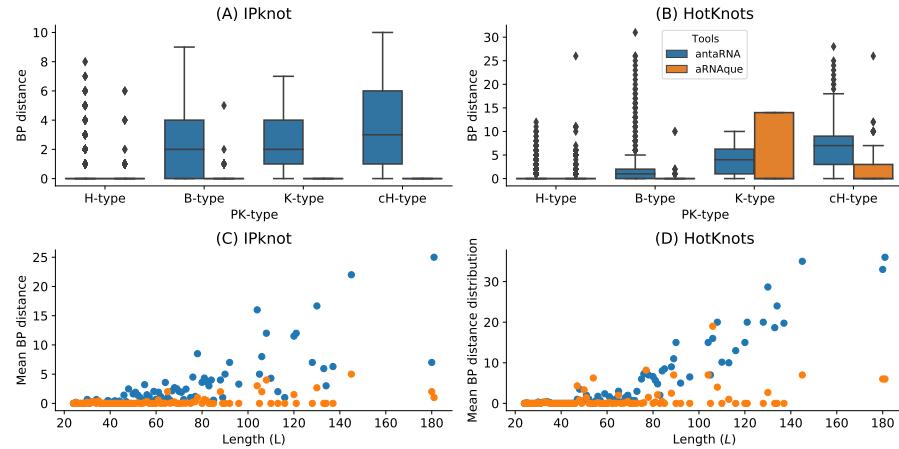


Figure 5.6: aRNAque vs antaRNA on PseudoBase++ dataset using both IPknot and HotKnots. Lower values imply better performance. (A, B) Base pair distance distributions of the designed sequences to the target structure for different pseudoknot types. (C,D) Mean base pair distance against target lengths.

3519 in terms of designed sequences quality (base pair distance dis-
 3520 tributions, with a t -value ≈ -1.04 and p -value ≈ 0.16) but more
 3521 significant in terms of the average number of generations needed
 3522 for successful matches to target structures (with a t -value ≈ -3.6
 3523 and p -value ≈ 0.0004). This result demonstrates a substantial
 3524 gain in computational time when using a Lévy mutation scheme
 3525 instead of a purely local mutation.

3526 5.2.2.3 Performance on PseudoBase++: aRNAque vs. antaRNA

3527 We also compared the sequences designed using aRNAque (with
 3528 the Lévy mutation scheme) to those produced by antaRNA. [Figure 5.6A](#) and [Figure 5.6C](#) show the base pair distance distribu-
 3529 tion for each category of pseudoknotted target structure and
 3530 the mean of the base pair distance plotted against the length of
 3531 the target secondary structures. For antaRNA, and when using
 3532 IPknot as a folding tool, finding sequences that fold into the tar-
 3533 get becomes increasingly difficult with pseudoknot complexity
 3534 (median base-pair distance distribution increases). On the other
 3535 hand, aRNAque's performance improves as pseudoknot complex-
 3536 ity increases (e.g. the mean base-distance decreases with the
 3537 pseudoknot complexity).

3538 A second benchmark using HotKnots as a folding tool was per-
 3539 formed on the same dataset. For both aRNAque and antaRNA, the

more complex the pseudoknot motifs, the worse is the tool performance (median of the base-pair distance distribution increases). **Figure 5.6B** and **Figure 5.6D** show the base pair distance distributions with respect to the pseudoknot motifs for both aRNAque and antaRNA. Even though both performances degrade as target length increases, aRNAque (Lévy flight evolutionary search) performance remains almost constant for all the target lengths greater than 60.

5.2.3 Quality of the designed RNA sequences

In addition to the successful rate analysis, we assessed the quality of the designed RNA sequences by analysing both GC-content and diversity of the pseudoknotted dataset using IPknot. This section presents the results obtained and a comparison to antaRNA designed sequences.

5.2.3.1 GC-content analysis of the designed sequences using IPknot

The GC-content of an RNA sequence S measures the concentration of G-C nucleotide in S and influences its stability and biological function. Therefore, the ability of an inverse folding tool to control the GC-content is of vital importance for designing functional RNA sequences. Both antaRNA and aRNAque allow to control the GC-content at different levels of the optimization process: aRNAque through the mutation parameters P_C and P_N ; antaRNA with the parameter $tGC \in [0, 1]$. In this section, we compare the performance of each tool for fixed GC-content values and analyse each tool's ability to control the GC-content. For each pseudoknotted target structure in the PseudoBase++ dataset, four different GC-content values {0.25, 0.5, 0.75, 1}, a poll of 20 sequences is designed using IPknot as folding tool. That results in 5320 designed sequences for each GC-content value and tool. The number of successes is the total number of sequences that fold exactly into the given target structure (i.e. the designed sequence folds into a structure at base-pair distance 0 from the target structure). **Figure 5.7** shows respectively the base pair distance distributions, the GC distance distributions and the number of successes for both aRNAque and antaRNA. The results show that the performance (in terms of success number) varies considerably with the GC-content values for both tools, and the best performance is obtained for both tools with a GC-content value of 0.5. When comparing the GC-content distance (i.e absolute

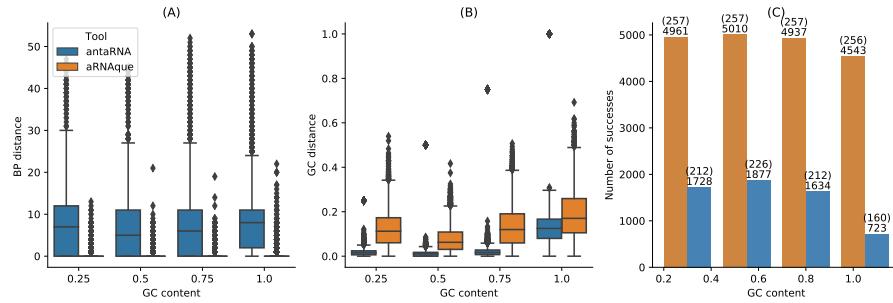


Figure 5.7: **aRNAque vs antaRNA on PseudoBase++ dataset using IPknot: GC-content analysis.** (A) Base-pair distance distributions. (B) GC-content distance distributions. The difference between the targeted GC-content and the actual GC-content values. In (A,B), lower values imply better performance. (C) Number of successes realised by both inverse folding tools. Two values are considered: the up value represent the number targets successfully solved for each GC-content value out of the 266 targets benchmarked; the down values represent the number sequences folding into the targeted secondary structure.

3580 value of the difference between the targeted GC–content and the
 3581 actual GC–content values of the designed sequences) distribu-
 3582 tions, both GC–content distance median distributions increase,
 3583 whereas antaRNA controls significantly better the GC–content
 3584 (See Figure 5.7B). On average, for the respective GC–content
 3585 values {0.25, 0.5, 0.75, 1}, antaRNA’s sequences have respectively
 3586 0.2569, 0.4952, 0.7314, 0.8684 whereas aRNAque’s sequences have
 3587 respectively 0.3649, 0.4910, 0.6231, 0.811; the main difference is
 3588 at fixed GC–content values 0.25 and 0.75. Even though antaRNA
 3589 designs sequences with better control of the GC–content, the gap
 3590 in success rate still remains remarkable compared to aRNAque (See
 3591 Figure 5.7A and Figure 5.7C).

3592 5.2.3.2 Diversity of the designed sequences

3593 Another advantage of using a Levy mutation when designing
 3594 RNA sequences is to increase the chance of designing sequences
 3595 with high diversity. Here, we use the positional entropy of each
 3596 pool of 20 sequences previously designed for each pseudoknot-
 3597 ted target structure to compare the diversity of RNA of both tools
 3598 antaRNA and aRNAque (Lévy search). We also compare it to the di-
 3599 versity of the designed sequences using the old version of aRNAque
 3600 (Local search). The results show that the sequence diversity of

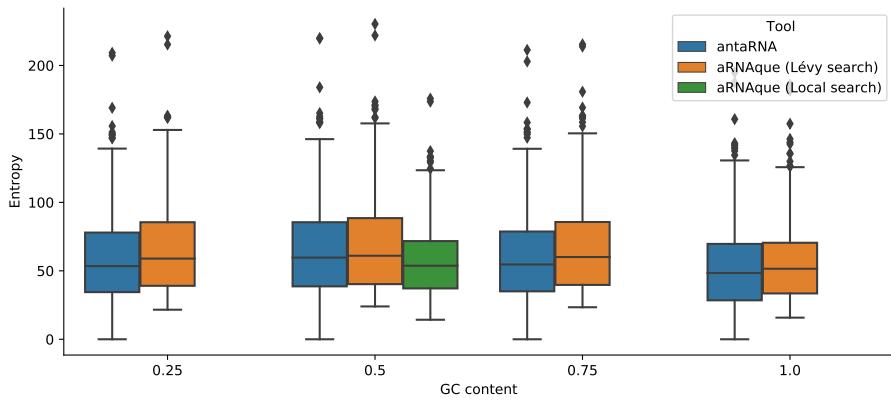


Figure 5.8: aRNAque vs antaRNA on PseudoBase++ dataset using IPknot: Diversity analysis. The positional entropy distributions plotted against the targeted GC-content values. Higher values imply better performance.

both antaRNA and aRNAque (Lévy search) varies with the GC-content values, where the more diversified pool of sequence is achieved with a GC-content value of 0.5. When comparing the pool of designed sequences with highest entropy (i.e. with a fixed GC-content of 0.5) to the one of the old version of aRNAque (Local search), the aRNAque (Lévy search) and antaRNA produce sequences with similar entropy (i.e. with a median entropy of 61.01 for Lévy search respectively 59.65 for antaRNA (see Figure 5.8), whereas the entropy of the sequences designed using the Local search is lower. For the three others fixed GC-content values (i.e. {0.25, 0.75, 1}, aRNAque (Lévy search) produces sequences with the highest entropy (respectively a median entropy of 58.9, 60.08, 51.52 against 53.42, 54.63, 48.38 for antaRNA).

5.2.4 Complexity and CPU time comparison

We finally analysed the design performance of aRNAque relatively to the CPU time needed. This section presents aRNAque statistical results compared to two main tools: RNAinverse for the pseudoknot-free targets and antaRNA for the pseudoknotted targets.

5.2.4.1 CPU time vs. success rate using RNAfold: RNAinverse vs. aRNAque on EteRNA100-V1.

Since our previous benchmarks on EteRNA100-V1 using the Turner2004 energy's parameters reveal that RNAinverse, one of the oldest in-

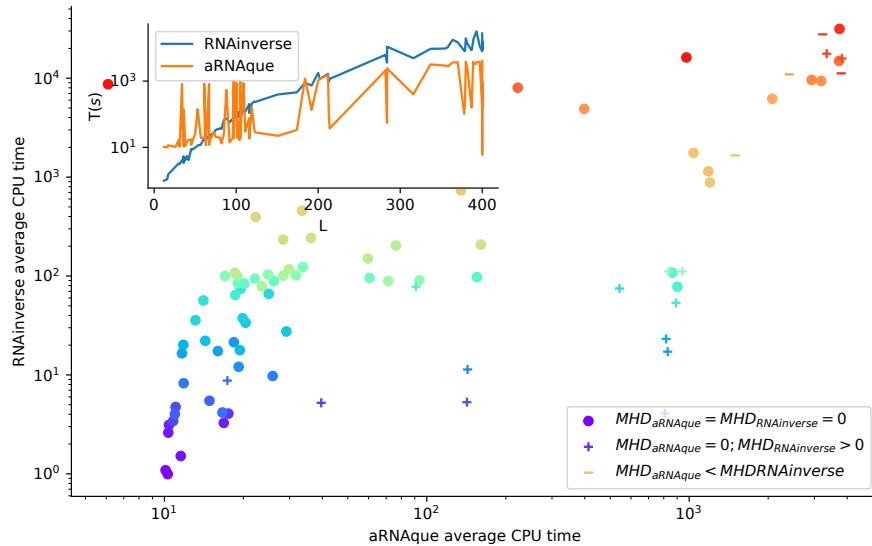


Figure 5.9: **CPU time: RNAinverse vs. aRNAque.** Each bubble corresponds to a target structure in EteRNA100 dataset and, their colours are proportional to the length of the targets. In the legend, MHD stands for Median Hamming distance, and the different markers represent—('o') 100% success for both tools—('+') 100% success for aRNAque and not for RNAinverse—('−') for the case both tools fail to find at least one sequence that folds into the target. Underlying the CPU time difference is the inside plot that shows the CPU time (in seconds) as a target length function.

3624 verse folding tools, stands behind aRNAque solving 66% of the
 3625 dataset; we have chosen to compare its computational time to our
 3626 implementation (See [Table 5.1](#)). The inset of [Figure 5.9](#) shows
 3627 the CPU time in seconds needed to design for each target in the
 3628 EteRNA100-V1, 5 sequences. As the RNAinverse time increases
 3629 exponentially with the length of the target, the aRNAque one does
 3630 not.

3631 When comparing the ratio between the success rate and CPU
 3632 time, aRNAque mostly succeeded in finding at least one sequence
 3633 that folds to the target with lower CPU time costs for average
 3634 target lengths. In contrast, RNAinverse accuracy is lower, and the
 3635 CPU time is expensive. The increase in CPU time may be because
 3636 of the use of the partition function as the objective function.

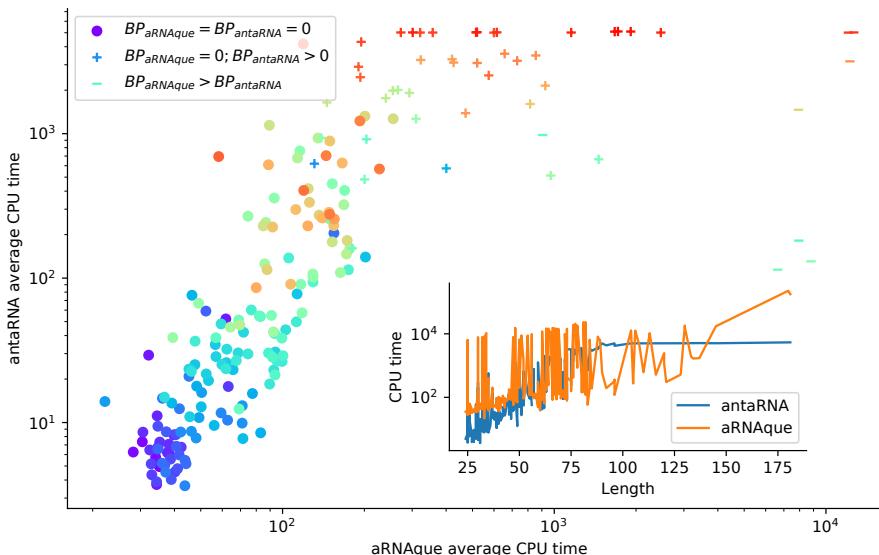


Figure 5.10: **CPU time analysis using Hotknots: antaRNA vs. aRNAque.**

Each bubble corresponds to a target structure in PseudoBase++ dataset and, their colours are proportional to the length of the targets. In the legend, BP stands for Median base pair distance, and the different markers represent—('o') 100% success for both tools—('+') 100% success for aRNAque and not for antaRNA—('−') for the case aRNAque’s desinged sequences are of median base pair distances greater than the one of antaRNA. Underlying the CPU time difference is the inside plot that shows the CPU time (in seconds) with respect to the target length.

3637 5.2.4.2 *CPU time vs. success rate using Hotknots: antaRNA vs.*
 3638 *aRNAque on PseudoBase++*

3639 We also compare aRNAque’s computational time to the one of
 3640 antaRNA. For both tools, 20 sequences were designed for each
 3641 target structure of the PseudoBase++ dataset. The GC-content
 3642 value used for both tools is 0.5, and the maximum number of
 3643 interactions for antaRNA is 5000. Figure 5.10 shows the median
 3644 CPU time of the 20 runs in seconds for both tools plotted against
 3645 each other. We analysed the CPU time by partitioning the data
 3646 into three groups: 1) a set for which both tools have a median
 3647 base-pair distance of 0 (158 entries marked with o); 2) another
 3648 set for which aRNAque has a median base-pair distance is 0 and
 3649 antaRNA (41 entries marked with +); 3) the last set for which
 3650 antaRNA designs are of better quality (9 entries mark as −). For
 3651 the first group, we can notice that for most targets of short length
 3652 antaRNA is faster than aRNAque. For the second group, although

3653 antaRNA average CPU time remains smaller, aRNAque's success rate
 3654 outperformed antaRNA. On the one hand, aRNAque average CPU
 3655 time is higher than the one of antaRNA, but this could be due to
 3656 its population-based algorithm, which often allows for designing
 3657 more successful sequences. On the other hand, antaRNA is faster
 3658 but less successful. Increasing antaRNA's number of iterations will
 3659 indeed increase the CPU time, but it may improve the quality of
 3660 the designed sequences.

3661 5.3 CONCLUSION

3662 In this work, we investigated an evolutionary approach to im-
 3663 prove the existing solutions to the RNA inverse folding problem.
 3664 As a result, we proposed a new EA python tool called aRNAque.
 3665 aRNAque implements a Lévy flight mutation scheme and supports
 3666 pseudoknotted RNA secondary structures. The benefit of a Lévy
 3667 flight over a purely local mutation search allowed us to explore
 3668 RNA sequence space at all scales. Such a heavy-tailed distribution
 3669 in the number of point mutations permitted the design of more
 3670 diversified sequences, avoiding the pitfalls of getting trapped in
 3671 a local optimum.

3672 Our results show general and significant improvements in the
 3673 design of RNA secondary structures compared to the standard
 3674 evolutionary algorithm mutation scheme with a mutation param-
 3675 eter $\approx 1/L$, where L is the sequence solution length. Lévy flight
 3676 mutations lead to a greater diversity of RNA sequence solutions
 3677 and, in many cases, reduce the evolutionary algorithm's number
 3678 of evaluations, thus improving computing time.

3679

Part III

3680

GENERAL CONCLUSION AND DISCUSSIONS

3681

3682

[August 4, 2022 at 19:02 – 1.0]

6

3683

3684 ADVANTAGES AND LIMITATIONS OF THE 3685 PROPOSED METHODS

3686 In the presented thesis, we have summarized some molecular
3687 background and biological functions of nucleic acids, especially
3688 ncRNAs. Because of the implication of the secondary structure
3689 of ncRNAs in performing biological functions and the separation
3690 of the folding time scale, our study focuses on the secondary
3691 structure of ncRNAs. Therefore, we have introduced the concepts
3692 of RNA bioinformatics and the essential computational problems
3693 related to the secondary structure of ncRNAs, such as RNA folding
3694 and the inverse folding. We presented a comprehensive literature
3695 review on existing tools that deal with both problems and some
3696 limitations for each tool. Despite advanced field results, we have
3697 introduced two new computation tools: RAFFT and aRNAque. What
3698 are the advantages and limitations of those tools? Is there any
3699 room for further improvements? How do these tools relate to
3700 evolutionary dynamics? In this concluding chapter, we will try
3701 to provide an answer to these questions by first discussing the
3702 advantages and the limitations of the tools previously introduced.

3703 6.1 rafft: LIMITATIONS AND FUTURE WORKS

3704 We presented in Chapter 3, RAFFT, a computational tool that ef-
3705 ficiently predicts RNA folding pathways. RAFFT takes advantage
3706 of the FFT to reduce its mean computational time to $O(N^2)$, es-
3707 pecially for long RNA sequences (length $\geq 10^3$). We assessed
3708 RAFFT performance for both the secondary structure prediction
3709 task and the RNA kinetics. In both cases, RAFFT shows important
3710 improvements. However, RAFFT also presents some limitations
3711 that will be addressed in the following section. We also discuss
3712 in this section some further improvements and applications.

3713 To first assess RAFFT performance for the folding task, two struc-
3714 ture estimates were compared with our method: the thermodynamic-
3715 based tools computed using RNAfold, LinearFold, RNAstructure
3716 and the ML estimate using MxFold2 and CONTRAfold. When we
3717 considered the lowest energy structure, the comparison of RAFFT
3718 to existing tools confirmed the overall validity of our approach.
3719 In more detail, a comparison with thermodynamic/ML models

3720 yielded the following results. First, the **ML** predictions performed
 3721 consistently better than both RAFFT and other approaches, where
 3722 the **PPV** = 70.4% and sensitivity = 77.1% on average. Second, the
 3723 **ML** methods produced loops, such as long hairpins or external
 3724 loops. We argue that the density of those loops correlates with
 3725 the ones in the benchmark dataset, which a **PCA** analysis revealed
 3726 too.

3727 In contrast, the density of similar loops was lower in the struc-
 3728 ture spaces produced by RAFFT and other thermodynamic-based
 3729 methods, implying some over-fitting in the **ML** model. Finally,
 3730 known structures obtained through covariation analysis reflect
 3731 *in vivo* structure conditions. Therefore, the structures predicted
 3732 by **ML** methods may result from their sequences alone and their
 3733 molecular environment, e.g. chaperones. We expect the thermo-
 3734 dynamic methods to provide a more robust framework for study-
 3735 ing sequence-to-structure relations. Concerning thermodynamic-
 3736 based tools, we obtained a substantial gain of performance when
 3737 analyzing $N = 50$ predicted structures per sequence, not only
 3738 the lowest energy one. This gain was even more remarkable for
 3739 sequences with fewer than 200 nucleotides, reaching the accuracy
 3740 of **ML** predictions.

3741 So how does RAFFT predictions contain structures that are more
 3742 relevant than the **MFE**, although these structures are less thermo-
 3743 dynamically stable? The interplay of three effects may explain
 3744 this finding. First, the **MFE** structure may not be relevant because
 3745 active structures can be in kinetic traps. Second, RAFFT forms a set
 3746 of pathways that cover the free energy landscape until they reach
 3747 local minima, yielding multiple long-lived structures accessible
 3748 from the unfolded state. Third, the energy function is not perfect,
 3749 so that the **MFE** structures computed by minimizing it may not in
 3750 fact be the most stable.

3751 We also showed that the fast-folding graph produced by RAFFT
 3752 can be used to reproduce state-of-the-art kinetics, at least qual-
 3753 itatively. Our method demonstrated three main benefits. First,
 3754 the kinetics can be drawn from as few as 68 structures, whereas
 3755 the barrier tree may require millions. Second, the kinetics ansatz
 3756 describes the complete folding mechanism starting from the un-
 3757 folded state. Third, for the length range tested here, the proce-
 3758 dure did not require any additional coarse-graining into basins.
 3759 (Longer **RNAs** might require such a coarse-graining step, in which
 3760 structures connected in the fast-folding graph are merged to-
 3761 gether).

3762 Based on our results, we believe that the proposed method is a
3763 robust heuristic for structure prediction and folding dynamics.
3764 The folding landscape depicted by RAFFT was designed to follow
3765 the kinetic partitioning mechanism, where multiple folding path-
3766 ways span the folding landscape. This approach has shown good
3767 predictive potential. Furthermore, we derived a kinetic ansatz
3768 from the fast-folding graph to model the slow part of the fold-
3769 ing dynamics. It was shown to approximate the usual kinetics
3770 framework qualitatively, although using many fewer structures.

3771 However, further improvements and extensions of the algo-
3772 rithm may be investigated. First, the choice of stems is limited
3773 to the largest in each positional lag, a greedy choice which may
3774 not be optimal. Second, we have constructed parallel pathways
3775 leading to diverse, accessible structures. Still, we have not given
3776 any thermodynamic-based criterion to identify which are more
3777 likely to resemble the native structure. We suggest using an **ML**-
3778 optimized score to this effect.

3779 Our method can also find applications in **RNA** design, where
3780 the design procedure could start with identifying long-lived in-
3781 termediates and using them as target structures. We also believe
3782 that mirror encoding can be helpful in phylogenetic analysis. In-
3783 deed, the correlation spectra $\text{cor}(k)$ computed here contained
3784 global information of base-pairing that can be used as a similarity
3785 measure.

3786 Finally, the versatile method implemented in RAFFT gives pos-
3787 sibilities for an alternative application of the **FFT** in **RNA-RNA**
3788 interaction. The underlying idea is that instead of encoding a
3789 sequence X and its mirror sequence \bar{X} , one can consider two
3790 encoded sequences X and Y , and the correlation between them
3791 will allow identifying the fraction of high interaction between
3792 two **RNA** sequences quickly. In general, **RNA-RNA** interaction pre-
3793 diction methods are divided into three groups: alignment like
3794 methods, **MFE** methods and comparative methods. **MFE** methods
3795 constitute the majority of the **RNA-RNA** interaction tools, with the
3796 only difference often based on whether the method considers
3797 intramolecular interactions. Some methods measure the acces-
3798 sibility of binding region (Intra and inter interactions) [7, 33,
3799 199]. We suggest neglecting intramolecular interactions and in-
3800 termolecular binding pairs for a preliminary implementation.

3801 In sum, RAFFT provides a versatile framework in which the
3802 kinetic partitioning mechanism can be simulated. Therefore, it
3803 allows for predicting an ensemble of concurrent **RNA** folding
3804 pathways ending in different metastable conformations. This

3805 result contrasts traditional thermodynamics techniques that find
 3806 a single MFE structure. However, further improvements of RAFFT
 3807 could be investigated:

- 3808 • the limitation of the choice of stems to the largest one in each
 3809 positional lag is a greedy choice that may not be optimal. We
 3810 propose to add stochastic noises in the choice of positional
 3811 lag to keep, such that running multiple times RAFFT, one
 3812 can overcome some greediness bottlenecks.
- 3813 • our method constructs parallel pathways leading to a di-
 3814 verse set of accessible structures. Still, we have not given
 3815 any thermodynamic-based criterion to identify which are
 3816 more likely to resemble the native structure. We suggest
 3817 using an ML-optimized score to investigate the restrained
 3818 ensemble of structures predicted by RAFFT.
- 3819 • structures connected in the parallel pathways are separated
 3820 by the formation or unfolding of a single stem. As men-
 3821 tioned above, RAFFT does not account for barriers between
 3822 structures that stem formation could involve. Therefore, we
 3823 propose to apply a post-treatment on the folding graph,
 3824 where the folding path between structures is investigated
 3825 using the set of valid atomic folding moves (*e.g.* individual
 3826 base-pair formation).

3827 In addition to these possible improvements, we presented two
 3828 possible applications: RNA design and RNA-RNA interactions. In
 3829 Section 6.3, we discuss another application in the study of evolu-
 3830 tionary dynamics.

3831 6.2 arnaque: LIMITATIONS AND PERSPECTIVES

3832 We have provided in Chapter 5, a new tool aRNAque, implementing
 3833 an EA with a Lévy flight mutation scheme that supports pseudo-
 3834 knotted RNA secondary structures. We discuss in this section the
 3835 advantages of using aRNAque for RNA design and some limitations
 3836 that could be addressed for further improvements.

3837 The Lévy mutation scheme offered exploration at different
 3838 scales (mostly local search combined with rare big jumps). Such a
 3839 scheme significantly improved the number of evaluations needed
 3840 to hit the target structure, while better avoiding getting trapped
 3841 in local optima. The benefit of a Lévy flight over a purely local
 3842 mutation search allowed us to explore RNA sequence space at all

scales. Such a heavy tailed distribution in the number of point mutations permitted the design of more diversified sequences. The main advantage of using a Lévy flight over local search was more remarkable for the pseudoknotted RNA targets, which is a reduction in the number of generations required to reach a target (see [Figure 5.5](#)). This is because the infrequent occurrence of a high number of mutations allow a diverse set of sequences among early generations, without the loss of robust local search. One consequence is a rapid increase in the population mean fitness over time and a rapid convergence to the target of the maximally fit sequence. To illustrate that advantage, we ran aRNAque starting from an initial population of unfolded sequences, both for a "one point mutation" and "Lévy mutation".

[Figure 6.1A](#) and [Figure 6.1B](#) show respectively the max/mean fitness over time and the number of distinct structures discovered over time plotted against the number of distinct sequences. When using a Lévy mutation scheme, the mean fitness increases faster in the beginning but stays lower than that using local mutations. Later in the optimisation, a big jump or high mutation on the RNA sequences produces structures with fewer similarities and, by consequence, worse fitness. In the (5–10)th generation, sequences folding into the target are already present in the Lévy flight population, but only at the 30th generation are similar sequences present in the local search population. The Lévy flight also allows exploration of both the structure and sequence spaces, providing a higher diversity of structures for any given set of sequences ([Figure 6.1B](#)). Using the mean entropy of structures as an alternate measure of diversity, we see in [Figure 6.1C](#) and [Figure 6.1D](#) how a Lévy flight achieves high diversity early in implementation, and maintains a higher diversity over all generations than a local search algorithm. Although the mutation parameters P_C and P_N influence the absolute diversity of the designed sequences, the Lévy flight always tends to achieve a higher relative diversity than local search, all else being equal.

We argue that the improved performance of the Lévy mutation over local search in target RNA structures is due to the high base pair density of pseudoknotted structures. Given that pseudoknotted RNA structures present a higher density of interactions, there are dramatic increases in possible incorrect folds and thus increasing risk of becoming trapped near local optima [[74](#)]. Large numbers of mutations in paired positions, as implied by a heavy tailed distribution, are necessary to explore radically different solutions.

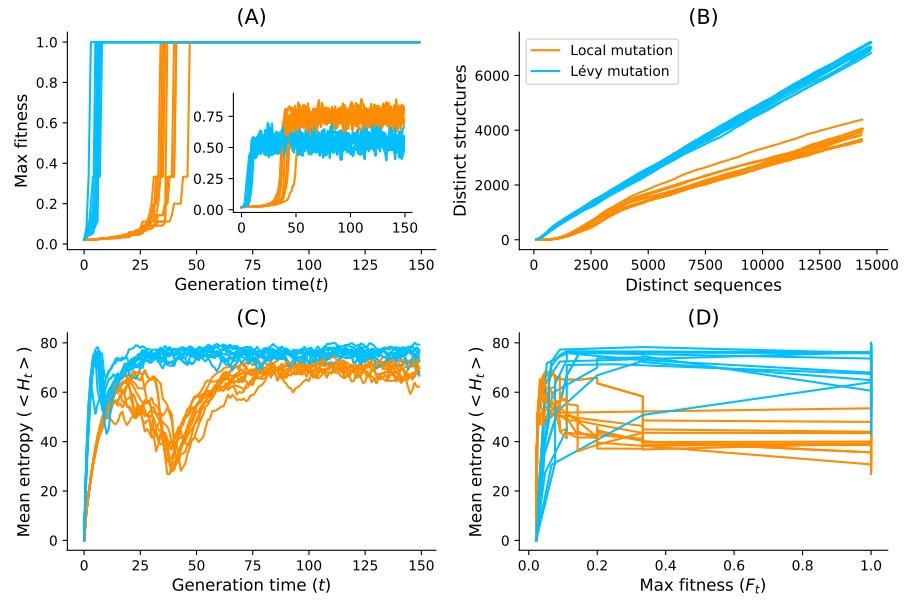


Figure 6.1: Lévy mutation vs one-point mutation. For the Eterna100 target structure [CloudBeta] 5 Adjacent Stack Multi-Branch Loop, ten independent runs were performed in which a minimum of 10 sequences were designed per run. (A) Max fitness and mean fitness (inset) over time. (B) Distinct sequences vs. Distinct structures over time. (C) Mean Shannon entropy of the population sequences over time for both binomial and Lévy mutation. (D) The max fitness plotted against the entropy over time.

3886 To illustrate that Lévy flight performance could be due to base
 3887 pair density, we clustered the benchmark datasets into two classes:
 3888 one cluster for target structures with low base pair density (den-
 3889 sity ≤ 0.5) and a second cluster for structures with high base pair
 3890 density (density > 0.5). Figure 5.3B showed the number of target
 3891 sequences available in each low and high density category. The
 3892 number of targets available in each category are colored according
 3893 to the percentage of pseudoknot-free targets (Eterna100-V1) vs.
 3894 targets with pseudoknots (Pseudobase++), showing that pseu-
 3895 doknots are strongly associated with high base pair densities:
 3896 71% of the pseudoknotted target structures have a high base pair
 3897 density. In contrast, the Eterna100 dataset without pseudoknots
 3898 has somewhat higher representation at low base pair density. If
 3899 it is true that improved Lévy flight performance is indeed tied to
 3900 base pair density, it is possible that similar heavy-tailed mutation
 3901 schemes could offer a scalable solution to even more complex
 3902 inverse folding problems. Another measure of difficulty is the
 3903 length of the target RNA secondary structure. When analysing the

3904 mean length of the pseudoknot-free targets, the high base-pair
3905 density targets are on average 181 nucleotides longer, and the low-
3906 density base-pair targets are 139 nucleotides (See [Figure 5.3C](#)).
3907 We have 49 nucleotides for low-density targets for the pseudo-
3908 knotted targets and 52 nucleotides for the high-density targets.
3909 That suggests that the Lévy mutation may be a good standard
3910 for designing more challenging target structures.

3911 A further effort have been made to understand the cases in
3912 which the Levy flight mutation can outperform the Binomial
3913 with low mutation rate or a constant one-point mutation rate.
3914 The key point of a Lévy mutation for the Inverse folding problem
3915 partially may rely on the base-pair density and the stability of
3916 stems with budge.

3917 Although we believe that Lévy flight-type search algorithms
3918 offer a valuable alternative to local search, we emphasise that its
3919 enhanced performance over say antaRNA is partially influenced by
3920 the specific capabilities of existing folding tools. Their limitations
3921 may account for the degradation of these tools as the pseudoknot
3922 motifs get increasingly complex (i.e. the incapacity of existing
3923 folding tools to predict some pseudoknot motifs influences the
3924 performance of both aRNAque and antaRNA). The Lévy mutation
3925 has also shown less potential in controlling the GC-content of
3926 the designed sequence when compared to antaRNA on pseudo-
3927 knotted target structures. antaRNA's parameters used in this work
3928 were tuned using pKiss; therefore, it could be possible room for
3929 improving the benchmark presented here by retuning them using
3930 IPKnot or HotKnots. Another possible limitation is the fact that
3931 most target structures were relatively easy to solve (in less than
3932 100 generations), which possibly allowed local search to perform
3933 better than Lévy search in some cases. Further research on more
3934 challenging target structures will improve our understanding of
3935 which conditions favour local *vs.* Lévy search.

3936 **6.3 rafft, arnaque AND EVOLUTIONARY DYNAMICS PERSPEC-**
3937 **TIVES**

3938 The [RNA](#) inverse folding has deep connections with theoretical
3939 evolutionary dynamics studies, where the sequence-secondary
3940 structure relationship is a popular model for studying the geno-
3941 type/phenotype maps [[66](#), [89](#)]. The folding tool usually maps
3942 each sequence to a secondary structure, e.g. RAFFT pathways
3943 could be used to compute developmental paths from sequence to
3944 secondary structure and then use the most dominant structure as

3945 the phenotype realization of a genotype RNA sequence. Therefore,
 3946 the two tools we previously introduced have a direct connection
 3947 with the evolutionary dynamic, where aRNAque simulates the
 3948 dynamic evolutionary process and RAFFT computes the genotype-
 3949 /phenotype mapping. This section presents some evolutionary
 3950 dynamics concepts that could be further study using RAFFT and
 3951 aRNAque.

3952 Similar to EAs, implemented in aRNAque, simulating a dynamic
 3953 evolutionary process using RNA sequence-secondary structure
 3954 relationship as a model often involves a population of RNA se-
 3955 quences to a given target secondary structure. In such a simula-
 3956 tion, we need three main ingredients: replication, selection and
 3957 mutation. These are the fundamental and defining principles of
 3958 biological systems. The underlying idea is that the genomic mate-
 3959 rial (the blueprint that determines the corresponding secondary
 3960 structure) in the form of RNAs is replicated and passed on to the
 3961 new offspring from generation to generation. An RNA individual
 3962 is then folded into its corresponding secondary structure at each
 3963 generation. Fitness is then a function that measures how close
 3964 is the realized structure to the target structure. Therefore, selec-
 3965 tion results when different types of RNA individuals compete
 3966 with each other. One RNA may reproduce faster and out-compete
 3967 the others. Occasionally, reproduction involves mistakes; these
 3968 mistakes are termed mutations. Mutations are responsible for
 3969 generating different RNAs that can be evaluated in the selection
 3970 process, thus resulting in biological novelty and diversity.

3971 Such a simple model gives a unified framework to precisely
 3972 define and statistically measure evolutionary dynamics concepts
 3973 such as plasticity, evolvability, epistasis, neutrality, continuity,
 3974 and modularity. At the molecular level, plasticity is viewed as the
 3975 capacity of an RNA sequence to assume a variety of energetically
 3976 favourable secondary structures by equilibrating among them
 3977 at a constant temperature [2]. Such concepts have been exten-
 3978 sively studied using the RNA inverse folding as a toy model. These
 3979 studies revealed that selection leads to the reduction of plasticity
 3980 and, therefore, to extreme modularity [2]. Another well-studied
 3981 property of evolution is neutrality which was first introduced
 3982 by Kimura [97], and it suggested that the majority of genotypic
 3983 changes (or mutations) in evolution are selectively neutral. The
 3984 attention to Kimura's contention has led to the discovery of neu-
 3985 tral networks in the context of genotype-phenotype models for
 3986 RNA secondary structure [145, 165]. Many recent studies [178,
 3987 179] use the sequence-secondary structure relationship as a toy

3988 model for studying neutral evolution. The neutral property of
3989 the RNA sequence-structure map contributes to a certain extent to
3990 the difficulty of the RNA design problem (e.g. when the neutral
3991 network is dense, this may quickly increase the chance of getting
3992 trapped and thus not improving the fitness). This problem is
3993 central to many optimization techniques and has already been
3994 mentioned in Chapter 4. Trying to avoid such a situation has
3995 motivated the choice of the mutation scheme implemented in
3996 aRNAque, which is the Lévy mutation.

3997 Another important issue in evolutionary biology concerns the
3998 extent to which the history of life has proceeded gradually or
3999 has been punctuated by discontinuous transitions at the level
4000 of phenotypes. Distinguishing the notion of continuous from
4001 discontinuous changes at the level of phenotypes requires a no-
4002 tion of nearness between phenotypes. This notion was previously
4003 introduced by Fontana and Peter [57], and it is based on the prob-
4004 ability of one phenotype being accessible from another through
4005 changes in the genotype. The RNA sequence-secondary structure
4006 relationship provides a framework where the notion of disconti-
4007 nuity transition is more precise. It allows understanding of how
4008 it arises in the model of evolutionary adaptation. This is done by
4009 simulating an RNA population that evolves toward a tRNA target
4010 secondary structure in a flow reactor logically constrained to a
4011 capacity of 1000 sequences. Once the secondary target structure
4012 is found, the evolutionary trajectory is backtracked to identify all
4013 the distinct structures involved and the transition between them.
4014 An example of continuous transition in Appendix (see Figure B.2)
4015 is the transition 18 → 10 whereas the transition 15 → 22 is said
4016 to be discontinuous.

4017 The simulation illustrated in Section B.8 was performed us-
4018 ing RNAfold, the folding tool included in the ViennaRNA package.
4019 When using ViennaRNA, the plastic ensemble of an RNA sequence
4020 ϕ is often considered to be the suboptimal ensemble structure Σ_ϕ
4021 within a user-defined energy range above the MFE at a constant
4022 temperature T . The ViennaRNA package provides an efficient tools
4023 RNAsubopt allowing to compute Σ_ϕ . In a more rigorous imple-
4024 mentation of plasticity, each of those structures in the ensemble
4025 Σ_ϕ should result from a developmental pathway. Therefore, the
4026 environmental changes may induce a change in the developmen-
4027 tal path, allowing switching from one structure in the structural
4028 ensemble to another. When considering the set of structures pro-
4029 duced using RAFFT, each meta-stable structure represents an RNA
4030 pathway; therefore, this ensemble can be considered a develop-

4031 mental plastic ensemble. Using RAFFT to simulate the evolution-
 4032 ary dynamic model may provide an alternative framework to
 4033 study evolutionary concepts like continuity and plasticity. Per-
 4034 haps, another way of defining continuous transition ($S_1 \rightarrow S_2$)
 4035 from structure S_1 to S_2 when using will be to check if the structure
 4036 S_2 is in the RAFFT's structure ensemble of the sequence with MFE
 4037 S_1 . In that wise, we suggest utilizing RAFFT to study and draw a
 4038 different interpretation of continuous evolutionary transition.

4039 **6.4 CONCLUSION**

4040 In sum, the two computational tools introduced in [Chapter 3](#) and
 4041 [Chapter 5](#) have been further explored. Both tools presented ad-
 4042 vantages and limitations, opening doors to further improvements
 4043 and applications.

4044 On the one hand, RAFFT predicts fast RNA pathways resulting
 4045 in an ensemble of metastable structures instead of a single MFE
 4046 structure implemented by most traditional methods. The ensem-
 4047 ble structures have the advantages of containing some structures
 4048 of biological relevance and reproducing complete kinetic sim-
 4049 lations of known RNAs. However, RAFFT method presents some
 4050 greediness in the choice of stems, does not provide any crite-
 4051 rion allowing to choose biological relevant structures from the
 4052 ensemble produced and does not account for barriers between
 4053 structures. Despite these limitations, RAFFT offers improvements
 4054 to the computational times and RNA kinetics, and its versatil-
 4055 ity opens the door to several applications from RNA design to
 4056 RNA-RNA interaction.

4057 On the other hand, aRNAque allowed designing RNA sequences
 4058 with higher diversity at a reduced number of evaluations for pseu-
 4059 doknotted target structures. Except for CPU time performance on
 4060 pseudoknotted targets, the success rate performance on both
 4061 pseudoknot-free and pseudoknotted target structures showed
 4062 improved performance. Despite different improvements, some
 4063 Eterna100 targets remain unsolvable, opening the door to fur-
 4064 ther investigations. We also discussed some aRNAque limitations,
 4065 such as the influence of the pseudoknot prediction capacities of
 4066 existing folding tools in the design process and aRNAque potential
 4067 to control the GC-content. In addition to these two limitations,
 4068 most pseudoknotted targets were solvable in less than 100 gener-
 4069 ations. These limitations contributed to the description of further
 4070 research directions.

⁴⁰⁷¹ Our results go beyond the computational RNA folding and
⁴⁰⁷² inverse folding; they can be used to study evolutionary dynamics
⁴⁰⁷³ concepts such as continuity and plasticity. Some perspectives
⁴⁰⁷⁴ have also been discussed.

[August 4, 2022 at 19:02 – 1.0]

4075

4076 GENERAL CONCLUSION

4077 This thesis has explored computational methods for studying
4078 RNA folding. In particular, it focused on the secondary structure
4079 level. It examined the energetic and thermodynamic stability
4080 characteristics in predicting folding pathways and designing RNA
4081 target structures through inverse folding. The principal output
4082 of the thesis was the development of computational tools to effi-
4083 ciently predict RNA folding pathways using the FFT (RAFFT) and
4084 an evolutionary algorithm allowing search at both local and long-
4085 range scales in the design of target RNA structures (aRNAque). On
4086 the one hand, our first contribution in RNA folding, RAFFT, offered
4087 an alternative computational framework to predict and study the
4088 RNA kinetics for long RNA molecules at lower computation costs
4089 than classical DP methods. The versatility of our methods opened
4090 doors to different ranges of applications, such as RNA-RNA inter-
4091 actions and evolutionary dynamics.

4092 On the other hand, our RNA inverse folding tool, aRNAque, of-
4093 fered a unified framework that combined the negative and posi-
4094 tive RNA design with an EA that implements a Lévy flight muta-
4095 tion scheme. Our results showed general and significant improve-
4096 ments in the design of RNA secondary structures (especially on
4097 the pseudoknotted targets) compared to the standard evolution-
4098 ary algorithm mutation scheme with a mutation parameter $\approx 1/L$,
4099 where L is the sequence solution length. Lévy flight mutation
4100 led to a greater diversity of RNA sequence solutions and reduced
4101 the evolutionary algorithm's number of evaluations, thus im-
4102 proving computing time compared to the local search. Although
4103 antaRNA average CPU time remains smaller, aRNAque's success
4104 rate outperformed antaRNA. To further improve our program, we
4105 suggest using a more powerful computational architecture such
4106 as massively parallel genetic algorithm (MPGA). This type of ar-
4107 chitecture may allow solving more challenging target secondary
4108 structures.

4109 Finally, we outlined these tools' limitations and prospects more
4110 generally in furthering our understanding of RNA structure, func-
4111 tion and design. We have put them into the context of evolution-
4112 ary dynamics and highlighted potential applications in studying
4113 continuous transitions and plasticity in that context. We believe

⁴¹¹⁴ that our contributions can enhance our understanding of RNA
⁴¹¹⁵ folding and could find applications in the real world—the exten-
⁴¹¹⁶ sions for future works and the implications for understanding
⁴¹¹⁷ their evolution.

4118

Part IV

4119

APPENDIX

4120

[August 4, 2022 at 19:02 – 1.0]

A

4121

4122 RAFFT APPENDICES

4123 A.1 KINETIC COMPARISON

4124 According to the RNA structure thermodynamics, one RNA molecule
4125 can adopt a structure s with probability $p(s) \propto \exp(-\beta\Delta G(s))$,
4126 where β is the inverse thermal energy (mol/kcal). To measure the
4127 quality of the ensemble of structures proposed by our method,
4128 we measured: (1) the average probability of each structures in
4129 the ensemble, then (2) the diversity of these structures.

4130 The probability coverage PC given by $PC(s) = \frac{1}{|\Omega|} \sum_{s \in \Sigma} p(s)$. Ω
4131 is the ensemble of structures sampled by a given method. We com-
4132 pared, for various random sequences, the probability coverage to
4133 methods based on Boltzmann sampling [46, 78]. We generated
4134 ensembles of 10^2 , 10^3 , and 10^4 structures per sequence denoted
4135 respectively SB100, SB1K, and SB10K. In addition, we also com-
4136 pared to RNAXplorer, a tool also based on a biased Boltzmann
4137 sampling.

4138 All structures are represented in the dot-bracket notation. In
4139 the dot-bracket notation, one structure has $ss = \{(.,.)\}$ symbols
4140 at each position. Given these three symbols, we propose the fol-
4141 lowing positional entropy measure $S = \frac{1}{L} \sum f_i(ss) \times \log(f_i(ss))$,
4142 where $f_i(ss)$ is the frequency of a symbol ss at position i in the
4143 ensemble of structure proposed.

4144 Figure A.1 shows the probability coverage and the positional
4145 entropy measure per method. It shows comparable sampling
4146 performances for fairly size sequences ($\approx 10^2$ nucleotides); and
4147 a comparable diversity.

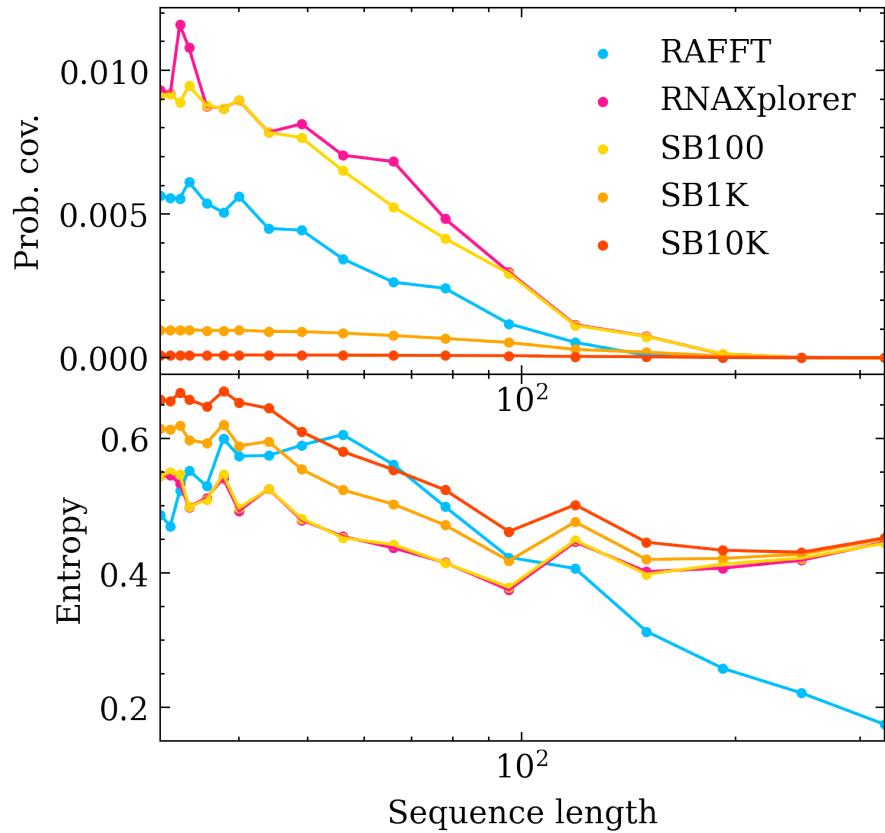


Figure A.1: Structure ensemble characterization. The upper part shows the average probability summed over the ensembles of structures predicted per sequence with different methods. The bottom part shows the average positional entropy of structures using the dot-bracket notation.

4148 A.2 rafft EXAMPLE CALLS

4149 RAFFT computes the fast-folding paths for a given sequence. Start-
 4150 ing from the completely unfolded structure, it quickly identifies
 4151 stems using the FFT-based technique.

4152 For the examples, in the Coronavirus frameshifting stimulation
 4153 element obtained from RFAM.

4154 To display only the final structures:

4155

Listing A.1: Command line to run RAFFT executable after installation

```
4156 $ rafft -s
4157     GGGUUUGCGGUGUAAGUGCAGCCGUCUUACACCGUGCGGCACAGGCACUAGUACUGAUGUCGUUAACAGGGCU
4158
4159 -ms 5
```

4160 _____
 4161 To display the visited/saved intermediates:
 4162 _____

Listing A.2: Command line to run RAFFT executable after installation

```
4163 $ rafft -s
4164   GGGUUUGCGGUGUAAGUGCAGCCGUCUUACACCGUGCGGCACAGGCACUAGUACUGAUGUCGUUAACAGGGCUUUUGACAU
4165
4166 -ms 5
4167 --traj
```

4168 _____
 4169 Here,
 4170 A result to this call could look like this:
 4171 _____

Listing A.3: RAFFT's output results

```
4172 GGGUUUGCGGUGUAAGUGCAGCCGUCUUACACCGUGCGGCACAGGCACUAGUACUGAUGUCGUUAACAGGGCUUUUGACAU
4173
4174 # -----0-----
4175 .....
4176     0.0
4177 # -----1-----
4178 .....(((((((((.....)))))))))) -----
4179 .....(.....) -14.0
4180 .....(.....) -----
4181 .....)))))..... -6.8
4182 .....(((((.....).....).....).....)
4183 .....)))))..... -6.4
4184 .....(((((.....).....).....).....)
4185 .....(.....) -5.5
4186 (((((.....).....).....))
4187 .....(.....)
4188 ..... -4.6
4189 # -----2-----
4190 .....((((((.....).....))))))(((((.....).....).....)))
4191 .....)))))..... -23.1
4192 .....((((((.....).....))))))(((((.....).....).....)))
4193 .....)))))..... -20.9
4194 .....(((((.....).....).....).....).....)
4195 .....))))).....))..... -18.8
4196 .....((((((.....).....))))))(((((.....).....).....)))
4197 .....)))))..... -18.7
4198 .....((((((.....).....))))))))))
4199 .....(((((.....).....).....).....)) -18.2
```

```

4200 # -----3-----
4201 .....((((((.((.....))))))))(((((((.((.....))))))
4202 )).... -24.0
4203 .....((((((.((.....))))))))(((((((.((.....))))))
4204 )).... -24.0
4205 .....((((((.((.....))))))))(((((((.((.....))))))
4206 )).... -23.1
4207 .....((((((.((.....))))))))(((((((.((.....))))))
4208 )).... -21.8
4209 .....((((((.((.....))))))))(((((((.((.....))))))
4210 )).... -20.9
4211
4212 GGGUUUGCGGUGUAAGUGCAGCCGUCUUACACCGUGCGGCACAGGCACUAGUACUGAUGUCGUAUACAGGGCUUUU
4213
4214 .....((((((.((.....))))))))(((((((.((.....))))))
4215 )).... -24.0
4216 .....((((((.((.....))))))))(((((((.((.....))))))
4217 )).... -24.0
4218 .....((((((.((.....))))))))(((((((.((.....))))))
4219 )).... -23.1
4220 .....((((((.((.....))))))))(((((((.((.....))))))
4221 )).... -21.8
4222 .....((((((.((.....))))))))(((((((.((.....))))))
4223 )).... -20.9
4224


---


4225 Where the columns shows respectively the predicted structures
4226 and their free energy.

```

4227 A.3 rafft PERFORMANCE ANALYSIS FOR 200 STRUCTURES
4228 SAVED.

4229 A.4 rafft PERFORMANCE ANALYSIS WITH VARIOUS VALUES
4230 OF MINIMUM ENERGY CONTRIBUTION REQUIRED FOR LOOP
4231 FORMATION

4232 A.5

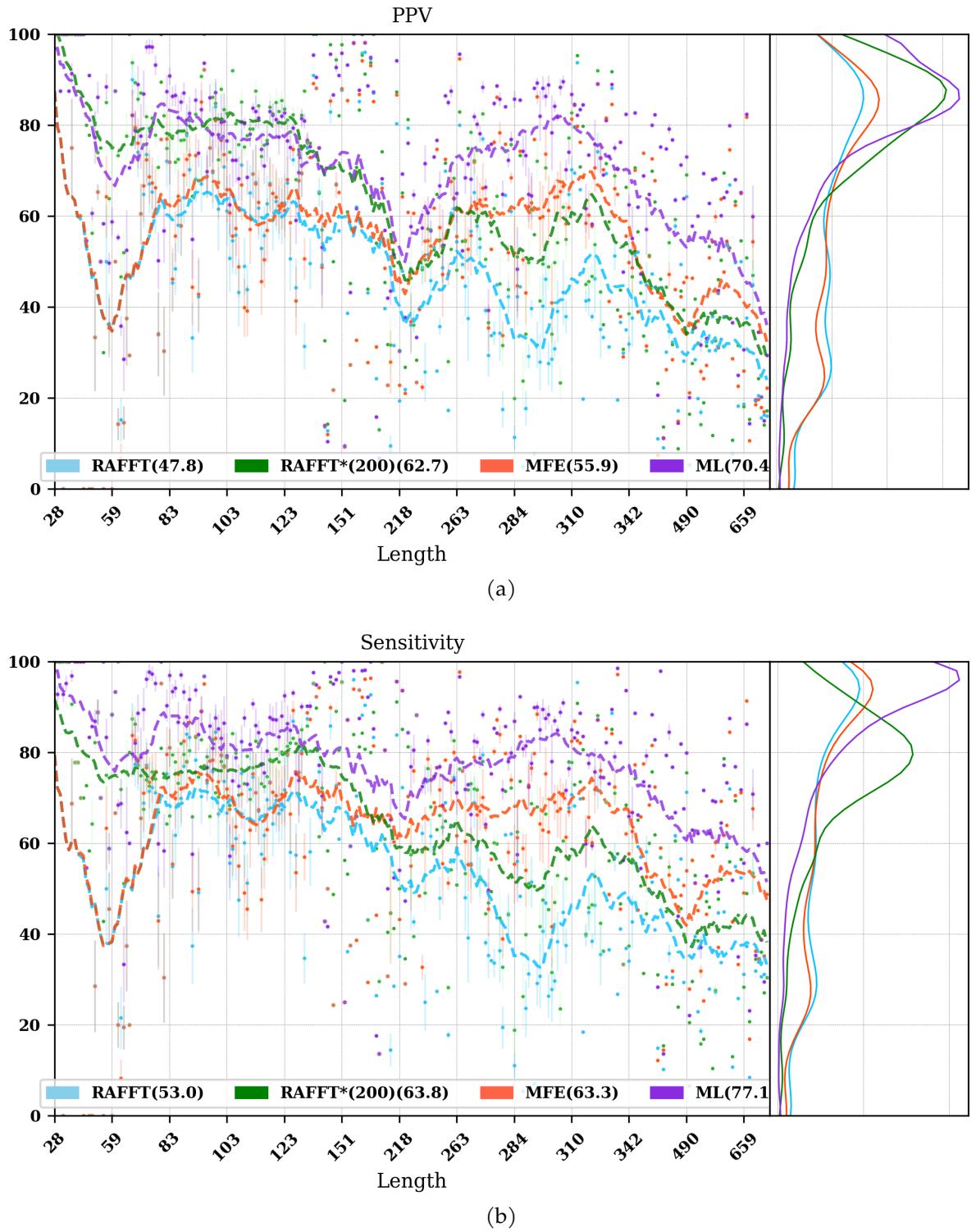


Figure A.2: Positive predictive values and sensitivity results. RAFFT (blue) displayed the best energy found. RAFFT*(200) shows the best score found among 200 saved structures. Left pans show the density (sequence-wise) of the accuracy measures.

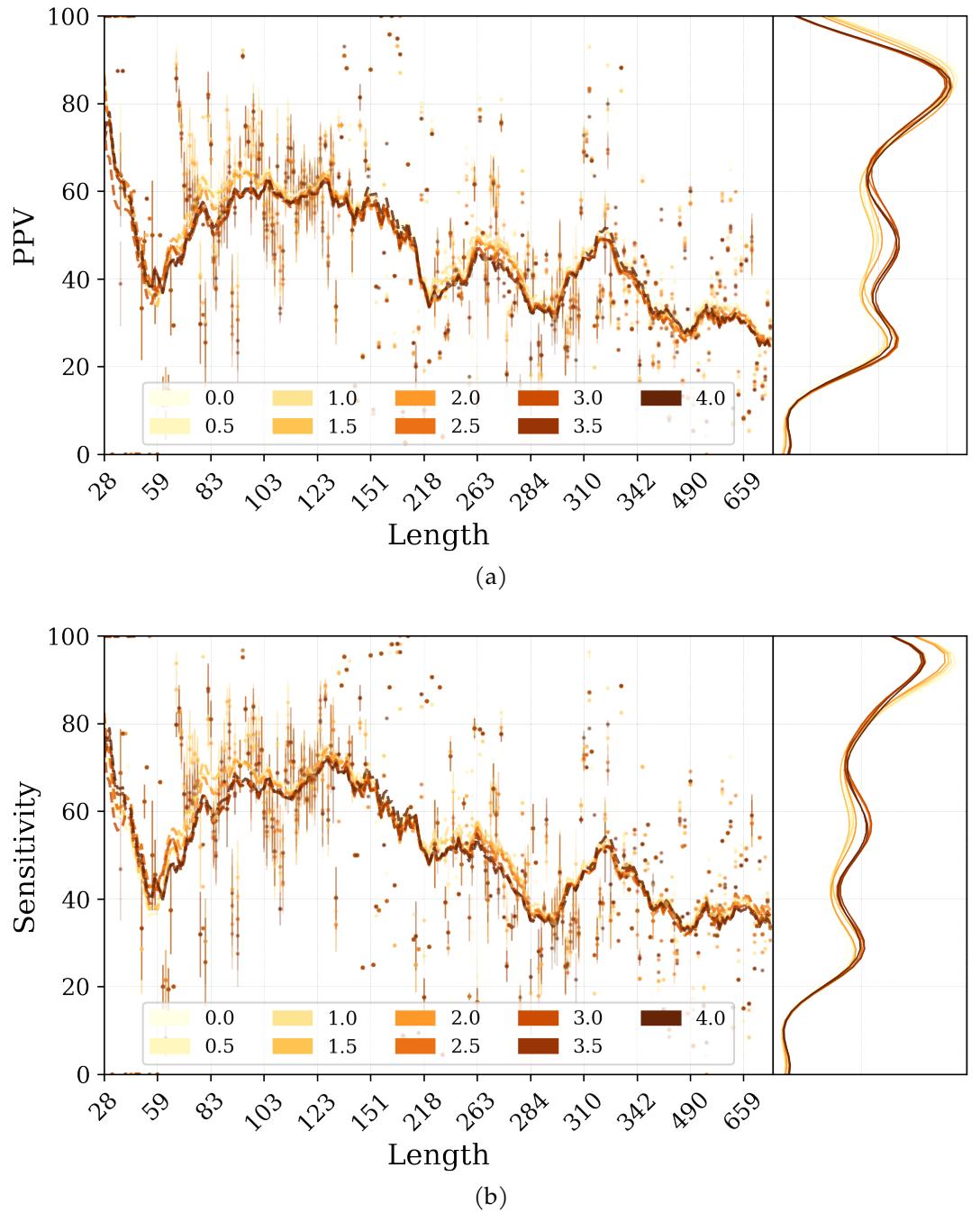


Figure A.3: Predictive performance of RAFFT with various values of minimum energy contribution required for loop formation. Positive values for this parameter causes RAFFT to accept destabilizing loops, therefore being less greedy than per default. The performance of RAFFT was not observed to be positively affected by allowing sub-optimal loop formation.

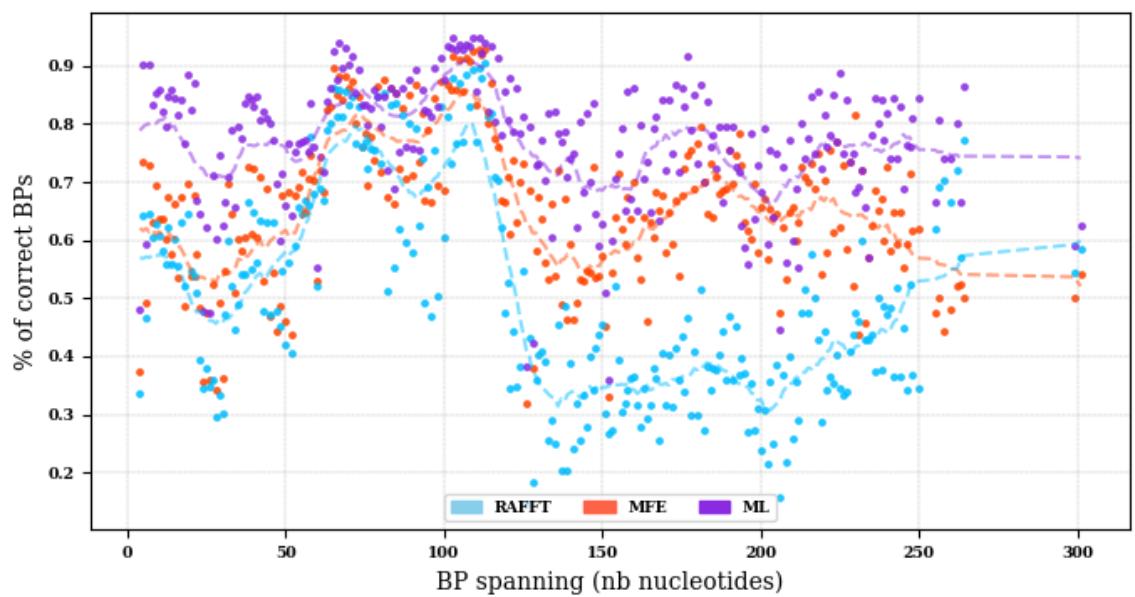


Figure A.4: **Base pair spanning:** It shows the percent of base pairs predicted found in the known structures per number of nucleotides between them.

[August 4, 2022 at 19:02 – 1.0]

B

4233

4234 ARNAQUE APPENDICES

4235 B.1 arnaque's gc-content parameters

4236 The GC-content is controlled in aRNAque using the mutation pa-
4237 rameters P_C and P_N . The following table gives the corresponding
4238 mutation parameters to the four regimes of GC-content values
4239 used for our benchmark.

Table B.1: Mutation parameters used in aRNAque to control the GC-content values.

GC-content values	P_C	P_N	aRNAque's key
0.25	{0.125, 0.125, 0.3, 0.3, 0.075, 0.075}	{0.125, 0.125, 0.375, 0.375}	GC25
0.25	{0.25, 0.25, 0.2, 0.2, 0.05, 0.05}	{0.25, 0.25, 0.25, 0.5}	GC50
0.75	{0.375, 0.375, 0.1, 0.1, 0.025, 0.025}	{0.375, 0.375, 0.125, 0.125}	GC75
1.0	{0.5, 0.5, 0.0, 0.0, 0.0, 0.0}	{0.5, 0.5, 0., 0.}	GC

4240 B.2 BENCHMARK ON eternal100 DATASET

4241 For each of the benchmarks on the Eterna100 datasets, We ran
4242 the first benchmark using the default aRNAque's parameter config-
4243 uration. And then, the unsolved structures are sorted out to run
4244 a second benchmark with a maximum number of generations
4245 set at 5000. aRNAque's performance presented in the paper is a
4246 combination of all the designed sequences for each realisation.

4247 B.3 GENERAL EA BENCHMARK PARAMETERS

4248 The same hardware resources and the same computer are used for
4249 all the benchmarks listed in the following table. A supercomputer
4250 with 40-Core Intel Xeon E5-2698 v4 at 2.2 GHz and 512 GB of
4251 RAM with a Debian OS.

Table B.2: Evolutionary algorithm parameter for each benchmarks.

Benchmark	Population size	# of generations (T)	Stopping criterion	Mutation parameter	# of runs per target
PseudoBase++ (IPknot)	100	200	$t = T$ $\max(f) = 0$	$c = 1.5$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	20
PseudoBase++ (Hotknots)	100	200	$t = T$ $\max(f) = 0$	$c = 1.5$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	20
PseudoBase++ GC-content (IPknot)	100	200	$t = T$ $\max(f) = 0$	$c = 1.5$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	20
Tuning Parameter (Binomial, IPknot)	100	200	$t = T$ $\max(f) = 0$	$\mu \in [0, 0.2]; c = None$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	20
Tuning Parameter (Lévy, IPknot)	100	200	$t = T$ $\max(f) = 0$	$c \in [1, 2]; \mu = None$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	20
Eterna100-V1 (OP, RNAfold)	100	5000	$t = T$ $\max(f) = 0$	$c = 7$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	5
Eterna100-V1 (Lévy, RNAfold)	100	5000	$t = T$ $\max(f) = 0$	$c = 1$ $P_N = \{0.7, 0.1, 0.1, 0.1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	5
Eterna100-V2 (OP, RNAfold)	100	5000	$t = T$ $\max(f) = 0$	$c = 7$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	5
Eterna100-V2 (Lévy, RNAfold)	100	5000	$t = T$ $\max(f) = 0$	$c = 1.5$ $P_N = \{0.7, 0.1, 0.1, .1\}$ $P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$	5

4252 B.4 OTHER BENCHMARK ON `eternal00-v1`

4253 The results on Eterna100-V1 presented in the paper are the best
 4254 of all the benchmarks we have performed. Since our mutation
 4255 scheme relies on the nucleotide distributions which implicitly
 4256 control the GC-content of the designed sequences, to obtain our
 4257 results, we first selected an arbitrary set of pairs $\{P_N, P_C\}$ and
 4258 benchmark aRNAque on Eterna100-V1 for each of them. The suc-
 4259 cess rate measures the fraction of sequences successfully folding
 4260 into the target structure. Table ?? shows the different parameters
 4261 we considered and the corresponded input key parameter using
 4262 the call of aRNAque script. Summary of the benchmark presented
 4263 in Table ?? is obtained by launching for each target structure 5
 4264 independent runs, with a population size of 100 and a maximum
 4265 number of generations of 5000. The energy parameter used here
 4266 was the Turner1999. The dashes in the table mean the benchmarks
 4267 have not been performed for the parameters.

Table B.3: Different parameters for the base pair distributions

Key	$P_N = \{p_A, p_G, p_U, p_C\}$	$P_C = \{p_{GC}, p_{CG}, p_{AU}, p_{UA}, p_{GU}, p_{UG}\}$
ALL	$P_N = \{0.25, 0.25, 0.25, 0.25\}$	$P_C = \{0.2, 0.2, 0.1, 0.1, 0.2, 0.2\}$
GC	$P_N = \{0.25, 0.25, 0.25, 0.25\}$	$P_C = \{0.5, 0.5, 0, 0, 0, 0\}$
GC ₁	$P_N = \{0.25, 0.65, 0.05, 0.05\}$	$P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$
GC ₂	$P_N = \{0.7, 0.1, 0.1, 0.1\}$	$P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$
GC ₃	$P_N = \{0.75, 0.1, 0.1, 0.05\}$	$P_C = \{0.4, 0.5, 0.1, 0, 0, 0\}$
GC ₄	$P_N = \{0.95, 0, 0.05, 0\}$	$P_C = \{0.4, 0.4, 0.2, 0, 0, 0\}$
GC ₅	$P_N = \{0.7, 0.1, 0.1, 0.1\}$	$P_C = \{0.3, 0.2, 0.2, 0.1, 0.1, 0.1\}$

Table B.4: Success percentage on Eterna100 datasets for each set of mutation parameters.

Tools	BP param	Mutation param	Percentage of success	$\#(Med(gen_{Zipf}) < Med(gen_{op}))$	$\#(Med(gen_{Zipf}) > Med(gen_{op}))$
aRNAque	<i>ALL</i>	Zipf ($c = 1$)	67%	7(#4)	
		One point	81%	64(#4)	
aRNAque	<i>GC</i>	Zipf ($c = 1$)	80%	43(#10)	
		One point	90%	30(#474)	
aRNAque	<i>GC₁</i>	Zipf ($c = 1$)	84%	29(#4)	
		One point	90%	33(#7)	
aRNAque	<i>GC₂</i>	Zipf ($c = 1$)	89%	61(#10)	
		One point	91%	19(#1920)	
aRNAque	<i>GC₃</i>	Zipf ($c = 1$)	88%	--	
		One point	--	--	
aRNAque	<i>GC₄</i>	Zipf ($c = 1$)	--	--	
		One point	--	--	
aRNAque	<i>GC₅</i>	Zipf ($c = 1$)	82%	44(#9)	
		One point	83%	30(#145)	
Total	–	Zipf ($c = 1$)	90%		
		One point	92%		
		RNAinverse	87%		

4268 B.5 TOOLS PATCHING

4269 To be able to perform our benchmarks, some slight modifications
 4270 was made on HotKnots and antaRNA. Details about the modifica-
 4271 tions are provided in this section.

- 4272 • **antaRNA:** The change was made at the line 1178 column 7,
 4273 where the line **args = 'HotKnots -m CC -s ' + sequence**
 4274 was replaced by to **args = './HotKnots -m CC -s ' + se-**
 4275 **quence.** The version of antaRNA we used is v2.0.1, and it
 4276 can be found on the Github link: <https://github.com/>
 4277 **RobertKleinkauf/antarna.**
- 4278 • **HotKnots:** to run HotKnots, we have to move aRNAque to the
 4279 bin directory. To avoid that, we updated the source code
 4280 and recompiled a new bin that does not require to move
 4281 aRNAque to the bin directory of HotKnots. We have uploaded

4282 the patched version of HotKnots in a third-part folder in
 4283 aRNAque's repository for benchmark reproduction.

4284 **NB:** The patches do not affect the folding algorithm. It con-
 4285 sisted of avoiding the use of relative paths in HotKnots.

4286 **B.6 arnaque EXAMPLE CALLS**

4287 aRNAque computes the RNA inverse folding problem for different
 4288 classes of structure complexities.

4289 For a pseudo-knot free target secondary structure:

Listing B.1: Command line to run aRNAque python script

```
4291 $ python aRNAque.py -t "((....)).((....))"  

  4292           -bp "GC2"  

  4293           -sm "NED"  

  4294           -ft "v"  

  4295           --job 5
```

4296 Here,

4297 A result to this call could look like this:

Listing B.2: aRNAque's output results

```
4300 GCUACGGCACCGUCAGG | ((....)).((....)) | -2.8 | 1.0  

  4301 GGGGGACCACCGGUGGG | ((....)).((....)) | -2.5 | 1.0  

  4302 GGGCCACCAGCGAAAGC | ((....)).((....)) | -2.2 | 1.0  

  4303 GGAAAUCCACCGGAAGG | ((....)).((....)) | -1.4 | 1.0  

  4304 GCAAGAGCGCCGCAAGG | ((....)).((....)) | -1.2 | 1.0
```

4305
 4306 Where the columns shows respectively the designed sequences,
 4307 the MFE structures, their free energy and the fitness to the target
 4308 (See Equation 5.1)

4309 **B.7 LÉVY FLIGHT VS LOCAL SEARCH: DESIGNING THE STRUCTURE
 4310 WITH THE SMALLEST NEUTRAL SET IN THE SPACE OF ALL
 4311 RNA SEQUENCES OF LENGTH 12**

4312 To further illustrate that advantage, we considered the space of
 4313 all RNA sequences of length 12 and with only G,C nucleotides.
 4314 The structures with the lowest neutral set are:

4315 1. $T_1 = (((...)).))$: only 2 sequences fold into the secondary
 4316 structure T_1

4317 2. $T_2 = ((.((...))))$: only 1 sequence folds into the secondary
 4318 structure T_2

4319 When having a close look at those two structures the base pair
 4320 density is maximal and there is an unpaired position on both that
 4321 allows the formation of a bulge.

4322 What that means naively is that any compatible sequence to T_1
 4323 (or T_2) will likely fold into a stem with four or three base pairs(
 4324 $((((...))).)$). Or $((((....))).)$, and these particular structures have
 4325 respectively 243 and 249 sequences in their neutral sets.

4326 We claim that, when having such kind of structure (T_1 or T_2),
 4327 the Levy mutation is of an important role to get out of the huge
 4328 neutral network of more stable stems. A simple test case was to
 4329 run aRNAque for a target secondary structure T_1 . For both one
 4330 point and Levy mutations, the distribution of the number of
 4331 generations needed to find sequences that fold into T_1 for both
 4332 mutation schemes is plotted in [Figure B.1](#).

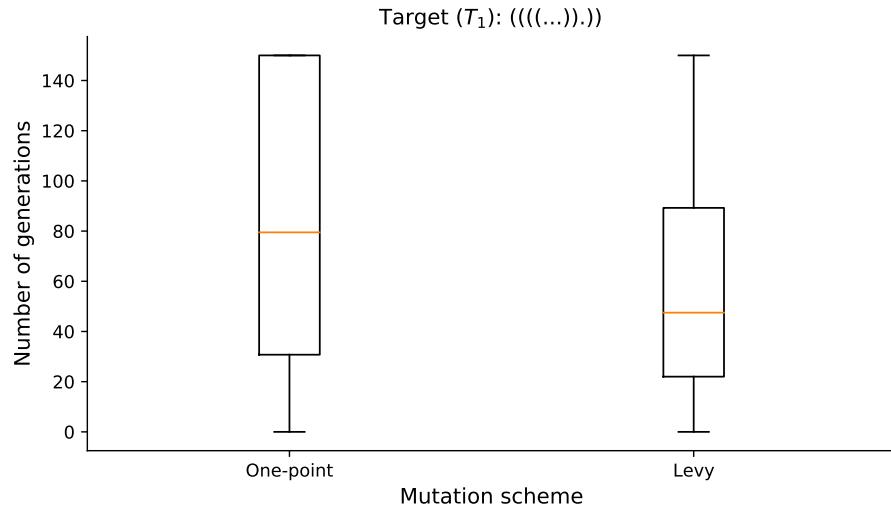
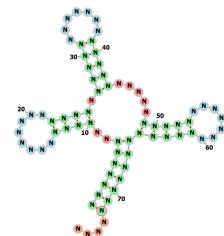


Figure B.1: Distribution of number of generations need to solve the target T_1 , for both Lévy and Local mutation schemes.

4333 B.8 CONTINUOUS AND DISCONTINUOUS TRANSITIONS IN EVO-
 4334 LUTION

4335 [Figure B.2](#) shows the evolution of the average distance to the
 4336 tRNA target structure, the intervals of time for which a particular
 4337 structure is present in the population, and a transition between

4338 distinct structures present in the evolutionary path. In Fontana's
4339 suggestions, a transition ($S_1 \rightarrow S_2$) between two structures S_1
4340 and S_2 is considered to be continuous if the structure S_1 is 'near'
4341 S_2 . In other terms, S_2 is likely to be accessible through the neigh-
4342 bor neutral sequences of S_1 . So if S_2 appears in the evolutionary
4343 path at time t , there exists a time $t' < t$ where S_2 was already
4344 present in the population. In contrast, the transition is discon-
4345 tinuous otherwise (i.e. the time the structure S_2 appears in the
4346 evolutionary path exactly at the same time it was present in the
4347 population). An example of continuous transition in [Figure B.2](#)
4348 is the transition $18 \rightarrow 10$ whereas the transition $15 \rightarrow 22$ is said
4349 to be discontinuous.



*tRNA target
secondary structure..*

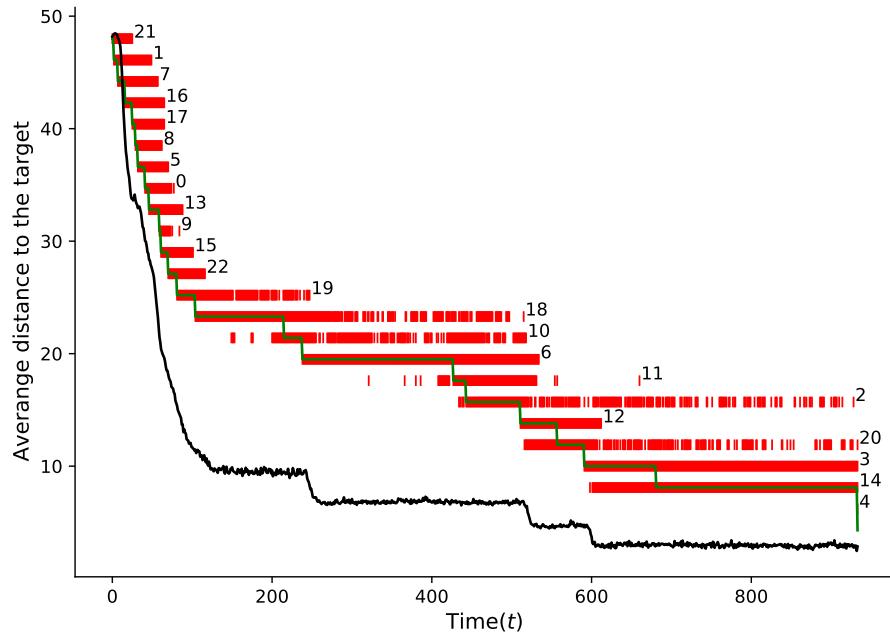


Figure B.2: Simulation of an RNA population evolving toward a tRNA (See the figure on the right side) target secondary structure. The target was reached after 933 generations (i.e. $\approx 10^5$ replications). The black line shows the average structure distance of the structures in the population to the target structure. The evolutionary history linking the initial structure to the target structure comprises 23. Each structure is labelled by an integer taken from 0 to 22. To each of them corresponds one horizontal line (in red). The top-level corresponds to the initial structure and the bottom the target structure. At each level, a series of red intervals correspond to the periods when the structure was present in the population, and the green curve represents the transition between structures. Only the time axis has a meaning for the red and green curves.

- 4351 [1] Fabian Amman, Stephan H. Bernhart, Gero Doose, Ivo L.
4352 Hofacker, Jing Qin, Peter F. Stadler, and Sebastian Will.
4353 “The trouble with long-range base pairs in RNA folding.”
4354 In: *Advances in Bioinformatics and Computational Biology*.
4355 Advances in Bioinformatics and Computational Biology.
4356 Springer International Publishing, 2013, pp. 1–11. doi:
4357 [10.1007/978-3-319-02624-4_1](https://doi.org/10.1007/978-3-319-02624-4_1). URL: https://doi.org/10.1007/978-3-319-02624-4_1.
- 4359 [2] Lauren W Ancel and Walter Fontana. “Plasticity, evolv-
4360 ability, and modularity in RNA.” In: *Journal of Experi-
4361 mental Zoology* 288.3 (2000), pp. 242–283.
- 4362 [3] Jeff Anderson-Lee, Eli Fisker, Vineet Kosaraju, Michelle
4363 Wu, Justin Kong, Jeehyung Lee, Minjae Lee, Mathew
4364 Zada, Adrien Treuille, and Rhiju Das. “Principles for pre-
4365 dicting RNA secondary structure design difficulty.” In: *Journal of molecular biology* 428.5 (2016), pp. 748–757.
- 4367 [4] Mirela Andronescu, Vera Bereg, Holger H Hoos, and
4368 Anne Condon. “RNA STRAND: the RNA secondary struc-
4369 ture and statistical analysis database.” In: *BMC Bioinfor-
4370 matics* 9.1 (2008), pp. 1–10.
- 4371 [5] Mirela Andronescu, Anthony P Fejes, Frank Hutter, Hol-
4372 ger H Hoos, and Anne Condon. “A new algorithm for
4373 RNA secondary structure design.” In: *Journal of molecular
4374 biology* 336.3 (2004), pp. 607–624.
- 4375 [6] Assaf Avihoo, Alexander Churkin, and Danny Barash.
4376 “RNAexinv: An extended inverse RNA folding from shape
4377 and physical attributes to sequences.” In: *BMC bioinfor-
4378 matics* 12.1 (2011), pp. 1–8.
- 4379 [7] Rolf Backofen and Wolfgang R Hess. “Computational
4380 prediction of sRNAs and their targets in bacteria.” In:
4381 *RNA biology* 7.1 (2010), pp. 33–42.
- 4382 [8] Rodolphe Barrangou, Christophe Fremaux, Hélène De-
4383 veau, Melissa Richards, Patrick Boyaval, Sylvain Moineau,
4384 Dennis A Romero, and Philippe Horvath. “CRISPR pro-
4385 vides acquired resistance against viruses in prokaryotes.”
4386 In: *Science* 315.5819 (2007), pp. 1709–1712.

- 4387 [9] S. Bellaousov and D. H. Mathews. "Probknot: fast pre-
 4388 diction of RNA secondary structure including pseudo-
 4389 knots." In: *RNA* 16.10 (2010), pp. 1870–1880. doi: 10 .
 4390 1261 / rna . 2125310. URL: <https://doi.org/10.1261/rna.2125310>.
- 4392 [10] Richard Bellman. "Dynamic programming, princeton
 4393 univ." In: *Press Princeton, New Jersey* (1957).
- 4394 [11] Jan H Bergmann and David L Spector. "Long non-coding
 4395 RNAs: modulators of nuclear structure and function." In:
 4396 *Current opinion in cell biology* 26 (2014), pp. 10–18.
- 4397 [12] Philip C Bevilacqua, Laura E Ritchey, Zhao Su, and Sarah
 4398 M Assmann. "Genome-wide analysis of RNA secondary
 4399 structure." In: *Annual review of genetics* 50 (2016), pp. 235–
 4400 266.
- 4401 [13] Eckart Bindewald, Kirill Afonin, Luc Jaeger, and Bruce A
 4402 Shapiro. "Multistrand RNA secondary structure predic-
 4403 tion and nanostructure design including pseudoknots."
 4404 In: *ACS nano* 5.12 (2011), pp. 9542–9551.
- 4405 [14] Édouard Bonnet, Paweł Rzążewski, and Florian Sikora.
 4406 "Designing RNA secondary structures is hard." In: *Journal*
 4407 of *Computational Biology* 27.3 (2020), pp. 302–316.
- 4408 [15] Ronald R Breaker, RF Gesteland, TR Cech, and JF Atkins.
 4409 *The RNA world*. Cold Spring Harbor Laboratory Press,
 4410 New York, 2006.
- 4411 [16] Philippe Brion and Eric Westhof. "Hierarchy and dynam-
 4412 ics of RNA folding." In: *Annual review of biophysics and*
 4413 *biomolecular structure* 26.1 (1997), pp. 113–137.
- 4414 [17] James W Brown. "The ribonuclease P database." In: *Nu-*
 4415 *cleic Acids Research* 26.1 (1998), pp. 351–352.
- 4416 [18] Anke Busch and Rolf Backofen. "INFO-RNA—a fast ap-
 4417 proach to inverse RNA folding." In: *Bioinformatics* 22.15
 4418 (2006), pp. 1823–1831.
- 4419 [19] Thomas R Cech and Joan A Steitz. "The noncoding RNA
 4420 revolution—trashing old rules to forge new ones." In: *Cell*
 4421 157.1 (2014), pp. 77–94.

- [20] Shaon Chakrabarti, Changbong Hyeon, Xiang Ye, George H Lorimer, and D Thirumalai. "Molecular chaperones maximize the native state yield on biological times by driving substrates out of equilibrium." In: *Proceedings of the National Academy of Sciences* 114.51 (2017), E10919–E10927.
- [21] James Chappell, Kyle E Watters, Melissa K Takahashi, and Julius B Lucks. "A renaissance in RNA synthetic biology: new mechanisms, applications and tools for the future." In: *Current opinion in chemical biology* 28 (2015), pp. 47–56.
- [22] Shi-Jie Chen. "RNA folding: conformational statistics, folding kinetics, and ion electrostatics." In: *Annu. Rev. Biophys.* 37 (2008), pp. 197–214.
- [23] Alexander Churkin, Matan Drory Retwitzer, Vladimir Reinhartz, Yann Ponty, Jérôme Waldspühl, and Danny Barash. "Design of RNAs: comparing programs for inverse RNA folding." In: *Briefings in bioinformatics* 19.2 (2017), pp. 350–358.
- [24] Simona Cocco, John F Marko, and Remi Monasson. "Slow nucleic acid unzipping kinetics from sequence-defined barriers." In: *The European Physical Journal E* 10.2 (2003), pp. 153–161.
- [25] James W Cooley and John W Tukey. "An algorithm for the machine calculation of complex Fourier series." In: *Mathematics of computation* 19.90 (1965), pp. 297–301.
- [26] Francis Crick. "Central dogma of molecular biology." In: *Nature* 227.5258 (1970), pp. 561–563.
- [27] Emilio Cusanelli and Pascal Chartrand. "Telomeric non-coding RNA: telomeric repeat-containing RNA in telomere biology." In: *Wiley Interdisciplinary Reviews: RNA* 5.3 (2014), pp. 407–419.
- [28] Paul Dallaire and François Major. "Exploring alternative RNA structure sets using MC-flashfold and db2cm." In: *RNA Structure Determination*. Springer, 2016, pp. 237–251.
- [29] Simon H Damberger and Robin R Gutell. "A comparative database of group I intron structures." In: *Nucleic Acids Research* 22.17 (1994), pp. 3508–3510.

- [4459] [30] Christian Darabos, Mario Giacobini, Ting Hu, and Jason H. Moore. "Lévy-Flight Genetic Programming: Towards a New Mutation Paradigm." In: *Evolutionary Computation, Machine Learning and Data Mining in Bioinformatics*. Ed. by Mario Giacobini, Leonardo Vanneschi, and William S. Bush. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 38–49. ISBN: 978-3-642-29066-4.
- [4460]
- [4461]
- [4462]
- [4463]
- [4464]
- [4465]
- [4466] [31] Kévin Darty, Alain Denise, and Yann Ponty. "VARNA: Interactive drawing and editing of the RNA secondary structure." In: *Bioinformatics* 25.15 (2009), p. 1974.
- [4467]
- [4468]
- [4469] [32] Jennifer Daub, Paul P Gardner, John Tate, Daniel Ram-sköld, Magnus Manske, William G Scott, Zasha Weinberg, Sam Griffiths-Jones, and Alex Bateman. "The RNA WikiProject: community annotation of RNA families." In: *RNA* 14.12 (2008), pp. 2462–2464.
- [4470]
- [4471]
- [4472]
- [4473]
- [4474] [33] Christoph Dieterich and Peter F Stadler. "Computational biology of RNA interactions." In: *Wiley Interdisciplinary Reviews: RNA* 4.1 (2013), pp. 107–120.
- [4475]
- [4476]
- [4477] [34] Ken A Dill. "Additivity principles in biochemistry." In: *Journal of Biological Chemistry* 272.2 (1997), pp. 701–704.
- [4478]
- [4479] [35] Robert M Dirks, Milo Lin, Erik Winfree, and Niles A Pierce. "Paradigms for computational nucleic acid de-sign." In: *Nucleic acids research* 32.4 (2004), pp. 1392–1403.
- [4480]
- [4481]
- [4482] [36] Robert M. Dirks and Niles A. Pierce. "A partition function algorithm for nucleic acid secondary structure includ-ing pseudoknots." In: *Journal of Computational Chemistry* 24.13 (2003). _eprint: <https://onlinelibrary.wiley.com/-doi/pdf/10.1002/jcc.10296>, pp. 1664–1677. ISSN: 1096-987X. DOI: <https://doi.org/10.1002/jcc.10296>. URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/jcc.10296> (visited on 02/17/2021).
- [4483]
- [4484]
- [4485]
- [4486]
- [4487]
- [4488]
- [4489]
- [4490] [37] Robert M Dirks and Niles A Pierce. "A partition function algorithm for nucleic acid secondary structure including pseudoknots." In: *Journal of computational chemistry* 24.13 (2003), pp. 1664–1677.
- [4491]
- [4492]
- [4493]
- [4494] [38] Chuong B Do, Daniel A Woods, and Serafim Batzoglou. "CONTRAFold: RNA secondary structure prediction with-out physics-based models." In: *Bioinformatics* 22.14 (2006), e90–e98.
- [4495]
- [4496]
- [4497]

- 4498 [39] Elizabeth A Doherty and Jennifer A Doudna. "Ribozyme
4499 structures and mechanisms." In: *Annual Review of Bio-*
4500 *physics and Biomolecular Structure* 30.1 (2001), pp. 457–
4501 475.
- 4502 [40] Ivan Dotu, Juan Antonio Garcia-Martin, Betty L Slinger,
4503 Vinodh Mechery, Michelle M Meyer, and Peter Clote.
4504 "Complete RNA inverse folding: computational design
4505 of functional hammerhead ribozymes." In: *Nucleic acids*
4506 *research* 42.18 (2014), pp. 11752–11762.
- 4507 [41] Robin D Dowell and Sean R Eddy. "Evaluation of several
4508 lightweight stochastic context-free grammars for RNA
4509 secondary structure prediction." In: *BMC bioinformatics*
4510 5.1 (2004), pp. 1–14.
- 4511 [42] D Draper, T Gluick, and P Schlax. In *RNA Structure and*
4512 *Function*, Vol. 298. 2000.
- 4513 [43] N Dromi, A Avihoo, and D Barash. "Reconstruction of nat-
4514 ural RNA sequences from RNA shape, thermodynamic
4515 stability, mutational robustness, and linguistic complexity
4516 by evolutionary computation." In: *Journal of Biomolecular*
4517 *Structure and Dynamics* 26.1 (2008), pp. 147–161.
- 4518 [44] Matan Drory Retwitzer, Vladimir Reinharz, Yann Ponty,
4519 Jérôme Waldispühl, and Danny Barash. "incaRNAbinv:
4520 a web server for the fragment-based design of RNA se-
4521 quences." In: *Nucleic acids research* 44.W1 (2016), W308–
4522 W314.
- 4523 [45] Andrew D Ellington, Xi Chen, Michael Robertson, and
4524 Angel Syrett. "Evolutionary origins and directed evolu-
4525 tion of RNA." In: *The international journal of biochemistry*
4526 & *cell biology* 41.2 (2009), pp. 254–265.
- 4527 [46] Gregor Entzian, Ivo L Hofacker, Yann Ponty, Ronny Lorenz,
4528 and Andrea Tanzer. "RNAXplorer: harnessing the power
4529 of guiding potentials to sample RNA landscapes." In:
4530 *Bioinformatics* 37.15 (2021), pp. 2126–2133.
- 4531 [47] Ali Esmaili-Taheri and Mohammad Ganjtabesh. "ERD:
4532 a fast and reliable tool for RNA design including con-
4533 straints." In: *BMC bioinformatics* 16.1 (2015), p. 20.
- 4534 [48] Ali Esmaili-Taheri, Mohammad Ganjtabesh, and Morteza
4535 Mohammad-Noori. "Evolutionary solution for the RNA
4536 design problem." In: *Bioinformatics* 30.9 (2014), pp. 1250–
4537 1258.

- 4538 [49] Manel Esteller. "Non-coding RNAs in human disease."
 4539 In: *Nature reviews genetics* 12.12 (2011), pp. 861–874.
- 4540 [50] Jörg Fallmann, Sebastian Will, Jan Engelhardt, Björn Grün-
 4541 ing, Rolf Backofen, and Peter F Stadler. "Recent advances
 4542 in RNA folding." In: *Journal of Biotechnology* 261 (2017),
 4543 pp. 97–104.
- 4544 [51] Alessandro Fatica and Irene Bozzoni. "Long non-coding
 4545 RNAs: new players in cell differentiation and develop-
 4546 ment." In: *Nature Reviews Genetics* 15.1 (2014), pp. 7–21.
- 4547 [52] Sven Findeiß, Manja Wachsmuth, Mario Mörl, and Peter F
 4548 Stadler. "Design of transcription regulating riboswitches."
 4549 In: *Methods in enzymology*. Vol. 550. Elsevier, 2015, pp. 1–
 4550 22.
- 4551 [53] Christoph Flamm, Walter Fontana, Ivo L Hofacker, and
 4552 Peter Schuster. "RNA folding at elementary step resolu-
 4553 tion." In: *Rna* 6.3 (2000), pp. 325–338.
- 4554 [54] Christoph Flamm, Ivo L Hofacker, Sebastian Maurer-
 4555 Stroh, Peter F Stadler, and Martin Zehl. "Design of multi-
 4556 stable RNA molecules." In: *Rna* 7.2 (2001), pp. 254–265.
- 4557 [55] Christoph Flamm, Ivo L. Hofacker, Peter F. Stadler, and
 4558 Michael T. Wolfinger. "Barrier trees of degenerate land-
 4559 scapes." In: *Zeitschrift für Physikalische Chemie* 216.2 (2002),
 4560 nil. doi: [10.1524/zpch.2002.216.2.155](https://doi.org/10.1524/zpch.2002.216.2.155). URL: <https://doi.org/10.1524/zpch.2002.216.2.155>.
- 4562 [56] Peter J Flor, James B Flanagan, and TR Cech. "A con-
 4563 served base pair within helix P4 of the Tetrahymena ri-
 4564 bozyme helps to form the tertiary structure required for
 4565 self-splicing." In: *The EMBO Journal* 8.11 (1989), pp. 3391–
 4566 3399.
- 4567 [57] Walter Fontana and Peter Schuster. "Continuity in evo-
 4568 lution: on the nature of transitions." In: *Science* 280.5368
 4569 (1998), pp. 1451–1455.
- 4570 [58] Jacques R. Fresco, Bruce M. Alberts, and Paul Doty. "Some
 4571 Molecular Details of the Secondary Structure of Ribonu-
 4572 cleic Acid." In: *Nature* 188.4745 (Oct. 1960). Number:
 4573 4745 Publisher: Nature Publishing Group, pp. 98–101.
 4574 ISSN: 1476-4687. doi: [10.1038/188098a0](https://doi.org/10.1038/188098a0). URL: <https://doi.org/10.1038/188098a0> (visited on
 4575 04/14/2021).
- 4576

- [59] James ZM Gao, Linda YM Li, and Christian M Reidys. "Inverse folding of RNA pseudoknot structures." In: *Algorithms for Molecular Biology* 5.1 (2010), pp. 1–19.
- [60] Juan Antonio Garcia-Martin, Peter Clote, and Ivan Dotu. "RNAiFOLD: a constraint programming algorithm for RNA inverse folding and molecular design." In: *Journal of bioinformatics and computational biology* 11.02 (2013), p. 1350001.
- [61] Paul P Gardner, Jennifer Daub, John G Tate, Eric P Nawrocki, Diana L Kolbe, Stinus Lindgreen, Adam C Wilkinson, Robert D Finn, Sam Griffiths-Jones, Sean R Eddy, et al. "Rfam: updates to the RNA families database." In: *Nucleic Acids Research* 37.suppl_1 (2009), pp. D136–D140.
- [62] Michael Geis, Christoph Flamm, Michael T Wolfinger, Andrea Tanzer, Ivo L Hofacker, Martin Middendorf, Christian Mandl, Peter F Stadler, and Caroline Thurner. "Folding kinetics of large RNAs." In: *Journal of Molecular Biology* 379.1 (2008), pp. 160–173.
- [63] David P Giedroc, Carla A Theimer, and Paul L Nixon. "Structure, stability and function of RNA pseudoknots involved in stimulating ribosomal frameshifting." In: *Journal of molecular biology* 298.2 (2000), pp. 167–185.
- [64] Sunny D Gilbert, Robert P Rambo, Daria Van Tyne, and Robert T Batey. "Structure of the SAM-II riboswitch bound to Sadenosylmethionine." In: *Nature structural & molecular biology* 15.2 (2008), pp. 177–182.
- [65] Walter Gilbert. "Origin of life: The RNA world." In: *Nature* 319.6055 (Feb. 1986). Number: 6055 Publisher: Nature Publishing Group, pp. 618–618. ISSN: 1476-4687. doi: [10.1038/319618a0](https://doi.org/10.1038/319618a0). URL: <https://www.nature.com/articles/319618a0> (visited on 04/09/2021).
- [66] Sam F Greenbury, Steffen Schaper, Sebastian E Ahnert, and Ard A Louis. "Genetic correlations greatly increase mutational robustness and can both reduce and enhance evolvability." In: *PLoS computational biology* 12.3 (2016), e1004773.
- [67] A. P. Gulyaev, F. H. van Batenburg, and C. W. Pleij. "An approximation of loop free energy values of RNA H-pseudoknots." In: *RNA (New York, N.Y.)* 5.5 (May 1999), pp. 609–617. ISSN: 1355-8382. doi: [10.1017/s135583829998189x](https://doi.org/10.1017/s135583829998189x).

- [68] Peixuan Guo. "The emerging field of RNA nanotechnology." In: *Nature nanotechnology* 5.12 (2010), pp. 833–842.
- [69] Zhuyan Guo and D Thirumalai. "Kinetics of protein folding: nucleation mechanism, time scales, and pathways." In: *Biopolymers: Original Research on Biomolecules* 36.1 (1995), pp. 83–102.
- [70] Robin R Gutell. "Collection of small subunit (16S-and 16S-like) ribosomal RNA structures: 1994." In: *Nucleic Acids Research* 22.17 (1994), pp. 3502–3507.
- [71] Robin R Gutell, Michael W Gray, and Murray N Schnare. "A compilation of large subunit (23S and 23S-like) ribosomal RNA structures: 1993." In: *Nucleic Acids Research* 21.13 (1993), p. 3055.
- [72] Robin R Gutell, Jung C Lee, and Jamie J Cannone. "The accuracy of ribosomal RNA comparative structure models." In: *Current opinion in structural biology* 12.3 (2002), pp. 301–310.
- [73] Robin R Gutell, Bryn Weiser, Carl R Woese, and Harry F Noller. "Comparative anatomy of 16-S-like ribosomal RNA." In: *Progress in nucleic acid research and molecular biology* 32 (1985), pp. 155–216.
- [74] Christine E Hajdin, Stanislav Bellaousov, Wayne Huggins, Christopher W Leonard, David H Mathews, and Kevin M Weeks. "Accurate SHAPE-directed RNA secondary structure modeling, including pseudoknots." In: *Proceedings of the National Academy of Sciences* 110.14 (2013), pp. 5498–5503.
- [75] Arif Ozgun Harmanci, Gaurav Sharma, and David H Mathews. "Stochastic sampling of the RNA structural alignment space." In: *Nucleic acids research* 37.12 (2009), pp. 4063–4075.
- [76] Christian Haslinger and Peter F. Stadler. "RNA structures with pseudo-knots: Graph-theoretical, combinatorial, and statistical properties." In: *Bulletin of Mathematical Biology* 61.3 (May 1, 1999), pp. 437–467. ISSN: 1522-9602. DOI: [10.1006/bulm.1998.0085](https://doi.org/10.1006/bulm.1998.0085). URL: <https://doi.org/10.1006/bulm.1998.0085> (visited on 04/15/2021).

- [77] Teresa Haynes, Debra Knisley, and Jeff Knisley. "Using a neural network to identify secondary RNA structures quantified by graphical invariants." In: *Comm Math Comput Chem* 60 (2008), pp. 277–290.
- [78] I. L. Hofacker. "Vienna RNA secondary structure server." In: *Nucleic Acids Research* 31.13 (2003), pp. 3429–3431. doi: [10.1093/nar/gkg599](https://doi.org/10.1093/nar/gkg599). URL: <https://doi.org/10.1093/nar/gkg599>.
- [79] Ivo L. Hofacker. "RNA Secondary Structure Prediction." In: *eLS*. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1038/npg.els.0005274>. American Cancer Society, 2005. ISBN: 978-0-470-01590-2. doi: [10.1038/npg.els.0005274](https://doi.org/10.1038/npg.els.0005274). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1038/npg.els.0005274> (visited on 03/03/2021).
- [80] Ivo L Hofacker, Walter Fontana, Peter F Stadler, L Sebastian Bonhoeffer, Manfred Tacker, and Peter Schuster. "Fast folding and comparison of RNA secondary structures." In: *Monatshefte für Chemie/Chemical Monthly* 125.2 (1994), pp. 167–188.
- [81] Ivo L. Hofacker, Peter F. Stadler, and Peter F. Stadler. "RNA Secondary Structures." In: *Reviews in Cell Biology and Molecular Medicine*. _eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/3527600906.mcb.200500009>. American Cancer Society, 2006. ISBN: 978-3-527-60090-8. doi: [10.1002/3527600906.mcb.200500009](https://doi.org/10.1002/3527600906.mcb.200500009). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1002/3527600906.mcb.200500009> (visited on 03/03/2021).
- [82] J Holland. *Adaptation in Natural and Artificial Systems*, Ann Arbor: University of Michigan Press, 1975. 1992.
- [83] Chiou-Yi Hor, Chang-Biau Yang, Chia-Hung Chang, Chiou-Ting Tseng, and Hung-Hsin Chen. "A Tool preference choice Method for RnA secondary structure prediction by SVM with statistical Tests." In: *Evolutionary Bioinformatics* 9 (2013), EBO-S10580.
- [84] Liang Huang, He Zhang, Dezhong Deng, Kai Zhao, Kaibo Liu, David A Hendrix, and David H Mathews. "LinearFold: linear-time approximate RNA folding by 5'-to-3' dynamic programming and beam search." In: *Bioinformatics* 35.14 (2019), pp. i295–i304.

- [85] Woong Y Hwang, Yanfang Fu, Deepak Reyon, Morgan L Maeder, Shengdar Q Tsai, Jeffry D Sander, Randall T Peterson, J-R Joanna Yeh, and J Keith Joung. "Efficient in vivo genome editing using RNA-guided nucleases." In: *Nature biotechnology* 31.3 (2013), p. 227.
- [86] Farren J Isaacs, Daniel J Dwyer, and James J Collins. "RNA synthetic biology." In: *Nature biotechnology* 24.5 (2006), pp. 545–554.
- [87] Hervé Isambert and Eric D Siggia. "Modeling RNA folding paths with pseudoknots: application to hepatitis delta virus ribozyme." In: *Proceedings of the National Academy of Sciences* 97.12 (2000), pp. 6515–6520.
- [88] Tor Ivry, Shahar Michal, Assaf Avihoo, Guillermo Sapiro, and Danny Barash. "An image processing approach to computing distances between RNA secondary structures dot plots." In: *Algorithms for Molecular Biology* 4.1 (2009), pp. 1–19.
- [89] Luc Jaeger, Eric Westhof, and Neocles B Leontis. "TectoRNA: modular assembly units for the construction of RNA nano-objects." In: *Nucleic acids research* 29.2 (2001), pp. 455–463.
- [90] Stefan Janssen and Robert Giegerich. "The RNA shapes studio." In: *Bioinformatics* (2015). doi: [10.1093/bioinformatics/btu649](https://doi.org/10.1093/bioinformatics/btu649). URL: <http://bioinformatics.oxfordjournals.org/content/31/3/423.abstract>.
- [91] Martin Jinek, Alexandra East, Aaron Cheng, Steven Lin, Enbo Ma, and Jennifer Doudna. "RNA-programmed genome editing in human cells." In: *elife* 2 (2013), e00471.
- [92] Steven G Johnson and Matteo Frigo. "A modified split-radix FFT with fewer arithmetic operations." In: *IEEE Transactions on Signal Processing* 55.1 (2006), pp. 111–119.
- [93] Anis Farhan Kamaruzaman, Azlan Mohd Zain, Suhaila Mohamed Yusuf, and Amirmudin Udin. "Lévy flight algorithm for optimization problems—a literature review." In: *Applied Mechanics and Materials*. Vol. 421. Trans Tech Publ. 2013, pp. 496–501.
- [94] Kazutaka Katoh, Kazuharu Misawa, Kei-ichi Kuma, and Takashi Miyata. "MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform." In: *Nucleic Acids Research* 30.14 (2002), pp. 3059–3066.

- [4733] [95] Michael Kertesz, Yue Wan, Elad Mazor, John L Rinn, Robert C Nutter, Howard Y Chang, and Eran Segal. "Genome-wide measurement of RNA secondary structure in yeast." In: *Nature* 467.7311 (2010), pp. 103–107.
- [4734]
- [4735]
- [4736]
- [4737] [96] Yoon Ki Kim, Luc Furic, Marc Parisien, François Major, Luc DesGroseillers, and Lynne E Maquat. "Staufen1 regulates diverse classes of mammalian transcripts." In: *The EMBO journal* 26.11 (2007), pp. 2670–2681.
- [4738]
- [4739]
- [4740]
- [4741] [97] Motoo Kimura. *The neutral theory of molecular evolution*. Cambridge University Press, 1983.
- [4742]
- [4743] [98] Daniel J Klein, Thomas E Edwards, and Adrian R Ferré-D'Amare. "Cocrystal structure of a class I preQ₁ riboswitch reveals a pseudoknot recognizing an essential hypermodified nucleobase." In: *Nature structural & molecular biology* 16.3 (2009), pp. 343–344.
- [4744]
- [4745]
- [4746]
- [4747]
- [4748] [99] Robert Kleinkauf, Torsten Houwaart, Rolf Backofen, and Martin Mann. "antaRNA—Multi-objective inverse folding of pseudoknot RNA using ant-colony optimization." In: *BMC bioinformatics* 16.1 (2015), pp. 1–7.
- [4749]
- [4750]
- [4751]
- [4752] [100] Robert Kleinkauf, Martin Mann, and Rolf Backofen. "antaRNA: ant colony-based RNA sequence design." In: *Bioinformatics* 31.19 (2015), pp. 3114–3121.
- [4753]
- [4754]
- [4755] [101] Konstantin Klemm, Christoph Flamm, and Peter F Stadler. "Funnels in energy landscapes." In: *The European Physical Journal B* 63.3 (2008), pp. 387–391.
- [4756]
- [4757]
- [4758] [102] Bjarne Knudsen and Jotun Hein. "RNA secondary structure prediction using stochastic context-free grammars and evolutionary history." In: *Bioinformatics (Oxford, England)* 15.6 (1999), pp. 446–454.
- [4759]
- [4760]
- [4761]
- [4762] [103] Bjarne Knudsen and Jotun Hein. "Pfold: RNA secondary structure prediction using stochastic context-free grammars." In: *Nucleic acids research* 31.13 (2003), pp. 3423–3428.
- [4763]
- [4764]
- [4765]
- [4766] [104] Donald E. Knuth. "Computer Programming as an Art." In: *Communications of the ACM* 17.12 (1974), pp. 667–673.
- [4767]
- [4768] [105] Rohan V Koodli, Benjamin Keep, Katherine R Coppess, Fernando Portela, Eterna participants, and Rhiju Das. "EternaBrain: Automated RNA design through move sets and strategies from an Internet-scale RNA videogame." In: *PLoS computational biology* 15.6 (2019), e1007059.
- [4769]
- [4770]
- [4771]
- [4772]

- [106] Rohan V. Koodli, Boris Rudolfs, Hannah K. Wayment-Steele, Eterna Structure Designers, and Rhiju Das. "Re-designing the EteRNA100 for the Vienna 2 folding engine." In: *bioRxiv* (2021). doi: [10.1101/2021.08.26.457839](https://doi.org/10.1101/2021.08.26.457839). eprint: [https://www.biorxiv.org/content/early/2021/08/28/2021.08.26.457839](https://www.biorxiv.org/content/early/2021/08/28/2021.08.26.457839.full.pdf). URL: <https://www.biorxiv.org/content/early/2021/08/28/2021.08.26.457839>.
- [107] Marcel Kucharík, Ivo L. Hofacker, Peter F. Stadler, and Jing Qin. "Pseudoknots in RNA folding landscapes." In: *Bioinformatics* 32.2 (Jan. 15, 2016). Publisher: Oxford Academic, pp. 187–194. ISSN: 1367-4803. doi: [10.1093/bioinformatics/btv572](https://doi.org/10.1093/bioinformatics/btv572). URL: <https://academic.oup.com/bioinformatics/article/32/2/187/1744549> (visited on 02/17/2021).
- [108] Yingjun Li, Saifu Pan, Yan Zhang, Min Ren, Mingxia Feng, Nan Peng, Lanming Chen, Yun Xiang Liang, and Qunxin She. "Harnessing Type I and Type III CRISPR-Cas systems for genome editing." In: *Nucleic acids research* 44.4 (2016), e34–e34.
- [109] Zhongsen Li, Zhan-Bin Liu, Aiqiu Xing, Bryan P Moon, Jessica P Koellhoffer, Lingxia Huang, R Timothy Ward, Elizabeth Clifton, S Carl Falco, and A Mark Cigan. "Cas9-guide RNA directed genome editing in soybean." In: *Plant physiology* 169.2 (2015), pp. 960–970.
- [110] Adam Lipowski and Dorota Lipowska. "Roulette-wheel selection via stochastic acceptance." In: *Physica A: Statistical Mechanics and its Applications* 391.6 (2012), pp. 2193–2196.
- [111] Qi Liu, Xiuzi Ye, and Yin Zhang. "A Hopfield neural network based algorithm for RNA secondary structure prediction." In: *First International Multi-Symposiums on Computer and Computational Sciences (IMSCCS'06)*. Vol. 1. IEEE. 2006, pp. 10–16.
- [112] Ronny Lorenz, Stephan H Bernhart, Christian Höner zu Siederdissen, Hakim Tafer, Christoph Flamm, Peter F Stadler, and Ivo L Hofacker. "Viennarna Package 2.0." In: *Algorithms for Molecular Biology* 6.1 (2011), p. 26. doi: [10.1186/1748-7188-6-26](https://doi.org/10.1186/1748-7188-6-26). URL: <https://doi.org/10.1186/1748-7188-6-26>.

- 4813 [113] Ronny Lorenz, Stephan H Bernhart, Christian Höner Zu
 4814 Siederdissen, Hakim Tafer, Christoph Flamm, Peter F
 4815 Stadler, and Ivo L Hofacker. "ViennaRNA Package 2.0."
 4816 In: *Algorithms for molecular biology* 6.1 (2011), p. 26.
- 4817 [114] Ronny Lorenz, Christoph Flamm, Ivo Hofacker, and Pe-
 4818 ter Stadler. "Efficient computation of base-pairing prob-
 4819 abilities in multi-strand RNA folding." In: *Proceedings*
 4820 of the 13th International Joint Conference on Biomedical En-
 4821 gineering Systems and Technologies. 2020, pp. 23–31. doi:
 4822 [10.5220/0008916600230031](https://doi.org/10.5220/0008916600230031). URL: <https://doi.org/10.5220/0008916600230031>.
- 4823 [115] Ronny Lorenz, Michael T Wolfinger, Andrea Tanzer, and
 4824 Ivo L Hofacker. "Predicting RNA secondary structures
 4825 from sequence and probing data." In: *Methods* 103 (2016),
 4826 pp. 86–98.
- 4827 [116] Zhi John Lu, Jason W Gloor, and David H Mathews. "Im-
 4828 proved RNA secondary structure prediction by maxi-
 4829 mizing expected pair accuracy." In: *Rna* 15.10 (2009),
 4830 pp. 1805–1813.
- 4831 [117] Rune B Lyngsø, James WJ Anderson, Elena Sizikova,
 4832 Amarendra Badugu, Tomas Hyland, and Jotun Hein. "Fr-
 4833 nakenstein: multiple target inverse RNA folding." In: *BMC*
 4834 *bioinformatics* 13.1 (2012), pp. 1–12.
- 4835 [118] Rune B. Lyngsø and Christian N. S. Pedersen. "RNA Pseu-
 4836 doknot Prediction in Energy-Based Models." In: *Journal of*
 4837 *Computational Biology* 7.3 (Aug. 1, 2000). Publisher: Mary
 4838 Ann Liebert, Inc., publishers, pp. 409–427. doi: [10.1089/106652700750050862](https://doi.org/10.1089/106652700750050862). URL: <https://doi.org/10.1089/106652700750050862> (visited on
 4841 04/16/2021).
- 4842 [119] JT Madison, GA Everett, and H Kung. "Nucleotide se-
 4843 quence of a yeast tyrosine transfer RNA." In: *Science* 153.3735
 4844 (1966), pp. 531–534.
- 4845 [120] B Mandelbrot. "Certain speculative prices (1963)." In:
 4846 *The Journal of Business* 45.4 (1972), pp. 542–543.
- 4847 [121] Hugo M. Martinez. "An RNA folding rule." In: *Nucleic*
 4848 *Acids Research* 12.1 (1984), pp. 323–334. doi: [10.1093/nar/12.1part1.323](https://doi.org/10.1093/nar/12.1part1.323). URL: <https://doi.org/10.1093/nar/12.1part1.323>.

- 4852 [122] David H. Mathews. "How to benchmark RNA secondary
 4853 structure prediction accuracy." In: *Methods* 162-163.162
 4854 (2019), pp. 60–67. doi: [10.1016/j.ymeth.2019.04.003](https://doi.org/10.1016/j.ymeth.2019.04.003).
 4855 URL: <https://doi.org/10.1016/j.ymeth.2019.04.003>.
- 4856 [123] David H Mathews, Matthew D Disney, Jessica L Childs,
 4857 Susan J Schroeder, Michael Zuker, and Douglas H Turner.
 4858 "Incorporating chemical modification constraints into a
 4859 dynamic programming algorithm for prediction of RNA
 4860 secondary structure." In: *Proceedings of the National Academy
 4861 of Sciences* 101.19 (2004), pp. 7287–7292.
- 4862 [124] David H Mathews, Jeffrey Sabina, Michael Zuker, and
 4863 Douglas H Turner. "Expanded sequence dependence of
 4864 thermodynamic parameters improves prediction of RNA
 4865 secondary structure." In: *Journal of Molecular Biology* 288.5
 4866 (1999), pp. 911–940.
- 4867 [125] DH Matthews, TC Andre, J Kim, DH Turner, and M Zuker.
 4868 "An updated recursive algorithm for RNA secondary
 4869 structure prediction with improved thermodynamic pa-
 4870 rameters." In: (1998).
- 4871 [126] Marco C Matthies, Stefan Bienert, and Andrew E Torda.
 4872 "Dynamics in sequence space for RNA secondary struc-
 4873 ture design." In: *Journal of chemical theory and computation*
 4874 8.10 (2012), pp. 3663–3670.
- 4875 [127] John S McCaskill. "The equilibrium partition function
 4876 and base pair binding probabilities for RNA secondary
 4877 structure." In: *Biopolymers: Original Research on Biomolecules*
 4878 29.6-7 (1990), pp. 1105–1119.
- 4879 [128] Nono SC Merleau and Matteo Smerlak. "A simple evo-
 4880 lutionary algorithm guided by local mutations for an
 4881 efficient RNA design." In: *Proceedings of the Genetic and
 4882 Evolutionary Computation Conference*. 2021, pp. 1027–1034.
- 4883 [129] Nono SC Merleau and Matteo Smerlak. "An evolution-
 4884 ary algorithm for inverse RNA folding inspired by Lévy
 4885 flights." In: *bioRxiv* (2022).
- 4886 [130] Gerard Minuesa, Cristina Alsina, Juan Antonio Garcia-
 4887 Martin, Juan Carlos Oliveros, and Ivan Dotu. "MoiR-
 4888 NAiFold: a novel tool for complex in silico RNA design."
 4889 In: *Nucleic acids research* 49.9 (2021), pp. 4934–4943.
- 4890 [131] Melanie Mitchell. *An introduction to genetic algorithms*. MIT
 4891 press, 1998.

- 4892 [132] Soheila Montaseri, Mohammad Ganjtabesh, and Fatemeh
 4893 Zare-Mirakabad. "Evolutionary algorithm for RNA sec-
 4894 ondary structure prediction based on simulated SHAPE
 4895 data." In: *PLoS one* 11.11 (2016), e0166965.
- 4896 [133] Peter B Moore and Thomas A Steitz. "The roles of RNA
 4897 in the synthesis of protein." In: *Cold Spring Harbor perspec-*
 4898 *tives in biology* 3.11 (2011), a003780.
- 4899 [134] Steffen Mueller, J Robert Coleman, Dimitris Papamichail,
 4900 Charles B Ward, Anjaruwee Nimnuan, Bruce Futcher,
 4901 Steven Skiena, and Eckard Wimmer. "Live attenuated
 4902 influenza virus vaccines by computer-aided rational de-
 4903 sign." In: *Nature biotechnology* 28.7 (2010), pp. 723–726.
- 4904 [135] Mark EJ Newman. "Power laws, Pareto distributions and
 4905 Zipf's law." In: *Contemporary physics* 46.5 (2005), pp. 323–
 4906 351.
- 4907 [136] Ruth Nussinov and Ann B Jacobson. "Fast algorithm for
 4908 predicting the secondary structure of single-stranded
 4909 RNA." In: *Proceedings of the National Academy of Sciences*
 4910 77.11 (1980), pp. 6309–6313.
- 4911 [137] Bibiana Onoa and Ignacio Tinoco Jr. "RNA folding and
 4912 unfolding." In: *Current opinion in structural biology* 14.3
 4913 (2004), pp. 374–379.
- 4914 [138] Vaitea Opuu, Nono SC Merleau, Messow Vincent, and
 4915 Matteo Smerlak. "RAFFT: Efficient prediction of RNA
 4916 folding pathways using the fast Fourier transform." In:
 4917 *bioRxiv* (2021).
- 4918 [139] Jie Pan, D. Thirumalai, and Sarah A. Woodson. "Folding
 4919 of RNA involves parallel pathways." In: *Journal of Molec-*
 4920 *ular Biology* 273.1 (1997), pp. 7–13. doi: [10.1006/jmbi.1997.1311](https://doi.org/10.1006/jmbi.1997.1311). URL: <https://doi.org/10.1006/jmbi.1997.1311>.
- 4923 [140] Marc Parisien and Francois Major. "The MC-Fold and
 4924 MC-Sym pipeline infers RNA structure from sequence
 4925 data." In: *Nature* 452.7183 (2008), pp. 51–55.
- 4926 [141] Fernando Portela. "An unexpectedly effective Monte Carlo
 4927 technique for the RNA inverse folding problem." 2018.
- 4928 [142] "RNACentral: a comprehensive database of non-coding
 4929 RNA sequences." In: *Nucleic acids research* 45.D1 (2017),
 4930 pp. D128–D134.

- 4931 [143] Effirul I Ramlan and Klaus-Peter Zauner. "Design of in-
 4932 teracting multi-stable nucleic acids for molecular infor-
 4933 mation processing." In: *Biosystems* 105.1 (2011), pp. 14–
 4934 24.
- 4935 [144] Jens Reeder, Peter Steffen, and Robert Giegerich. "pknot-
 4936 sRG: RNA pseudoknot folding including near-optimal
 4937 structures and sliding windows." In: *Nucleic acids research*
 4938 35.suppl_2 (2007), W320–W324.
- 4939 [145] Christian Reidys, Peter F Stadler, and Peter Schuster. "Generic
 4940 properties of combinatory maps: neutral networks of
 4941 RNA secondary structures." In: *Bulletin of mathematical
 4942 biology* 59.2 (1997), pp. 339–397.
- 4943 [146] Vladimir Reinharz, Yann Ponty, and Jérôme Waldspühl.
 4944 "A weighted sampling algorithm for the design of RNA
 4945 sequences with targeted secondary structure and nu-
 4946 cleotide distribution." In: *Bioinformatics* 29.13 (2013), pp. i308–
 4947 i315. ISSN: 1367-4803. DOI: [10.1093/bioinformatics/btt217](https://doi.org/10.1093/bioinformatics/btt217). eprint: <https://academic.oup.com/bioinformatics/article-pdf/29/13/i308/18534655/btt217.pdf>. URL:
 4948 <https://doi.org/10.1093/bioinformatics/btt217>.
- 4951 [147] Jihong Ren, Baharak Rastegari, Anne Condon, and Hol-
 4952 ger H Hoos. "HotKnots: heuristic prediction of RNA sec-
 4953 ondary structures including pseudoknots." In: *RNA* 11.10
 4954 (2005), pp. 1494–1504.
- 4955 [148] Jessica S Reuter and David H Mathews. "RNAsstructure:
 4956 software for RNA secondary structure prediction and
 4957 analysis." In: *BMC bioinformatics* 11.1 (2010), pp. 1–9.
- 4958 [149] Andy M Reynolds. "Current status and future directions
 4959 of Lévy walk research." In: *Biology open* 7.1 (2018), bio030106.
- 4960 [150] Elena Rivas and Sean R. Eddy. "A dynamic pro-
 4961 gramming algorithm for RNA structure prediction includ-
 4962 ing pseudoknots." In: *Journal of Molecular Biology* 285.5
 4963 (1999), pp. 2053–2068. ISSN: 0022-2836. DOI: <https://doi.org/10.1006/jmbi.1998.2436>. URL: <https://www.sciencedirect.com/science/article/pii/S0022283698924366>.
- 4966 [151] Elena Rivas, Raymond Lang, and Sean R Eddy. "A range
 4967 of complex probabilistic models for RNA secondary struc-
 4968 ture prediction that includes the nearest-neighbor model
 4969 and more." In: *RNA* 18.2 (2012), pp. 193–212.

- 4970 [152] Elena Rivas, Raymond Lang, and Sean R Eddy. "A range
 4971 of complex probabilistic models for RNA secondary struc-
 4972 ture prediction that includes the nearest-neighbor model
 4973 and more." In: *RNA* 18.2 (2012), pp. 193–212.
- 4974 [153] Debra L Robertson and Gerald F Joyce. "Selection in
 4975 vitro of an RNA enzyme that specifically cleaves single-
 4976 stranded DNA." In: *Nature* 344.6265 (1990), pp. 467–468.
- 4977 [154] John M Rosenberg, Nadrian C Seeman, Roberta O Day,
 4978 and Alexander Rich. "RNA double-helical fragments at
 4979 atomic resolution: II. The crystal structure of sodium
 4980 guanylyl-3', 5'-cytidine nonahydrate." In: *Journal of molec-*
 4981 *ular biology* 104.1 (1976), pp. 145–167.
- 4982 [155] Frederic Runge, Danny Stoll, Stefan Falkner, and Frank
 4983 Hutter. "Learning to design RNA." In: *arXiv preprint arXiv:1812.11951*
 4984 (2018).
- 4985 [156] Rick Russell, Xiaowei Zhuang, Hazen P Babcock, Ian S
 4986 Millett, Sebastian Doniach, Steven Chu, and Daniel Her-
 4987 schlag. "Exploring the folding landscape of a structured
 4988 RNA." In: *Proceedings of the National Academy of Sciences*
 4989 99.1 (2002), pp. 155–160.
- 4990 [157] Yasubumi Sakakibara, Michael Brown, Richard Hughey,
 4991 I Saira Mian, Kimmen Sjölander, Rebecca C Underwood,
 4992 and David Haussler. "Stochastic context-free grammars
 4993 for tRNA modeling." In: *Nucleic acids research* 22.23 (1994),
 4994 pp. 5112–5120.
- 4995 [158] Tore Samuelsson and Christian Zwieb. "The signal recog-
 4996 nition particle database (SRPDB)." In: *Nucleic Acids Re-*
 4997 *search* 27.1 (1999), pp. 169–170.
- 4998 [159] Baby Santosh, Akhil Varshney, and Pramod Kumar Ya-
 4999 dava. "Non-coding RNAs: biological functions and ap-
 5000 plications." In: *Cell biochemistry and function* 33.1 (2015),
 5001 pp. 14–22.
- 5002 [160] Kengo Sato, Manato Akiyama, and Yasubumi Sakakibara.
 5003 "RNA secondary structure prediction using deep learning
 5004 with thermodynamic integration." In: *Nature communica-*
 5005 *tions* 12.1 (2021), pp. 1–9.
- 5006 [161] Kengo Sato, Manato Akiyama, and Yasubumi Sakakibara.
 5007 *RNA secondary structure prediction using deep learning with*
 5008 *thermodynamic integration.* 2021.

- 5009 [162] Kengo Sato, Yuki Kato, Michiaki Hamada, Tatsuya Akutsu,
 5010 and Kiyoshi Asai. "IPknot: fast and accurate prediction
 5011 of RNA secondary structures with pseudoknots using
 5012 integer programming." In: *Bioinformatics* 27.13 (2011),
 5013 pp. i85–i93.
- 5014 [163] Martin Sauvageau, Loyal A Goff, Simona Lodato, Boyan
 5015 Bonev, Abigail F Groff, Chiara Gerhardinger, Diana B
 5016 Sanchez-Gomez, Ezgi Hacisuleyman, Eric Li, Matthew
 5017 Spence, et al. "Multiple knockout mouse models reveal
 5018 lincRNAs are required for life and brain development."
 5019 In: *elife* 2 (2013), e01749.
- 5020 [164] Murray N Schnare, Simon H Damberger, Michael W Gray,
 5021 and Robin R Gutell. "Comprehensive comparison of struc-
 5022 tural characteristics in eukaryotic cytoplasmic large sub-
 5023 unit (23 S-like) ribosomal RNA." In: *Journal of Molecular
 5024 Biology* 256.4 (1996), pp. 701–719.
- 5025 [165] Peter Schuster, Walter Fontana, Peter F Stadler, and Ivo L
 5026 Hofacker. "From sequences to shapes and back: a case
 5027 study in RNA secondary structures." In: *Proceedings of the
 5028 Royal Society of London. Series B: Biological Sciences* 255.1344
 5029 (1994), pp. 279–284.
- 5030 [166] Nadrian C Seeman, John M Rosenberg, FL Suddath, Jung
 5031 Ja Park Kim, and Alexander Rich. "RNA double-helical
 5032 fragments at atomic resolution: I. The crystal and molec-
 5033 ular structure of sodium adenylyl-3', 5'-uridine hexahy-
 5034 drate." In: *Journal of molecular biology* 104.1 (1976), pp. 109–
 5035 144.
- 5036 [167] Matthew G Seetin and David H Mathews. "RNA struc-
 5037 ture prediction: an overview of methods." In: *Bacterial
 5038 Regulatory RNA* (2012), pp. 99–122.
- 5039 [168] Martin J Serra and Douglas H Turner. "Predicting thermo-
 5040 dynamic properties of RNA." In: *Methods in enzymology*.
 5041 Vol. 259. Elsevier, 1995, pp. 242–261.
- 5042 [169] Bruce A Shapiro and Kaizhong Zhang. "Comparing mul-
 5043 tiple RNA secondary structures using tree comparisons."
 5044 In: *Bioinformatics* 6.4 (1990), pp. 309–318.
- 5045 [170] Vishnu Prakash Sharma, Harji Ram Choudhary, Sandeep
 5046 Kumar, and Vikas Choudhary. "A modified DE: Popula-
 5047 tion or generation based levy flight differential evolution
 5048 (PGLFDE)." In: *2015 International Conference on Futuristic*

- 5049 Trends on Computational Analysis and Knowledge Management (ABLAZE). IEEE. 2015, pp. 704–710.
- 5050
- 5051 [171] Jade Shi, Rhiju Das, and Vijay S Pande. “SentRNA: Im-
5052 proving computational RNA design by incorporating a
5053 prior of human design strategies.” 2018.
- 5054 [172] Micheal F Shlesinger, George M Zaslavsky, and Uriel
5055 Frisch. *Lévy flights and related topics in physics*. Vol. 450.
5056 Berlin, Heidelberg: Springer Berlin Heidelberg, 1995.
- 5057 [173] Wenjie Shu, Ming Liu, Hebing Chen, Xiaochen Bo, and
5058 Shengqi Wang. “ARDesigner: a web-based system for
5059 allosteric RNA design.” In: *Journal of biotechnology* 150.4
5060 (2010), pp. 466–473.
- 5061 [174] Christian Höner zu Siederdissen, Stephan H Bernhart, Pe-
5062 ter F Stadler, and Ivo L Hofacker. “A folding algorithm for
5063 extended RNA secondary structures.” In: *Bioinformatics*
5064 27.13 (2011), pp. i129–i136.
- 5065 [175] Jaswinder Singh, Jack Hanson, Kuldip Paliwal, and Yaoqi
5066 Zhou. “RNA secondary structure prediction using an
5067 ensemble of two-dimensional deep neural networks and
5068 transfer learning.” In: *Nature Communications* 10.1 (2019),
5069 pp. 1–13.
- 5070 [176] Michael F Sloma and David H Mathews. “Exact calcula-
5071 tion of loop formation probability identifies folding mo-
5072 tifs in RNA secondary structures.” In: *RNA* 22.12 (2016),
5073 pp. 1808–1818.
- 5074 [177] Michael F Sloma and David H Mathews. “Base pair prob-
5075 ability estimates improve the prediction accuracy of RNA
5076 non-canonical base pairs.” In: *PLoS computational biology*
5077 13.11 (2017), e1005827.
- 5078 [178] Matteo Smerlak. “Effective potential reveals evolution-
5079 ary trajectories in complex fitness landscapes.” In: *arXiv
5080 preprint arXiv:1912.05890* (2019).
- 5081 [179] Matteo Smerlak. “Neutral quasispecies evolution and the
5082 maximal entropy random walk.” In: *Science advances* 7.16
5083 (2021), eabb2376.
- 5084 [180] Sergey V. Solomatin, Max Greenfeld, Steven Chu, and
5085 Daniel Herschlag. “Multiple native states reveal persis-
5086 tent ruggedness of an RNA folding landscape.” In: *Nature*
5087 463.7281 (2010), pp. 681–684. doi: [10.1038/nature08717](https://doi.org/10.1038/nature08717).
5088 URL: <https://doi.org/10.1038/nature08717>.

- 5089 [181] T. Specht, M. Szymanski, M. Z. Barciszewska, J. Barciszewski,
 5090 and V. A. Erdmann. "Compilation of 5s rRNA and 5s
 5091 rRNA gene sequences." In: *Nucleic Acids Research* 25.1
 5092 (1997), pp. 96–97. doi: [10.1093/nar/25.1.96](https://doi.org/10.1093/nar/25.1.96). URL: <https://doi.org/10.1093/nar/25.1.96>.
- 5094 [182] Robert C Spitale, Andrew T Torelli, Jolanta Krucinska,
 5095 Vahe Bandarian, and Joseph E Wedekind. "The Structural
 5096 Basis for Recognition of the PreQo Metabolite by an Un-
 5097 usually Small Riboswitch Aptamer Domain." In: *Journal*
 5098 of Biological Chemistry 284.17 (2009), pp. 11012–11016.
- 5099 [183] Mathias Sprinzl, Carsten Horn, Melissa Brown, Anatoli
 5100 Ioudovitch, and Sergey Steinberg. "Compilation of tRNA
 5101 sequences and sequences of tRNA genes." In: *Nucleic*
 5102 *Acids Research* 26.1 (1998), pp. 148–153.
- 5103 [184] Evan W Steeg. "Neural networks, adaptive optimization,
 5104 and RNA secondary structure prediction." In: *Artificial*
 5105 *intelligence and molecular biology* (1993), pp. 121–60.
- 5106 [185] Paul R Stein and Michael S Waterman. "On some new se-
 5107 quences generalizing the Catalan and Motzkin numbers."
 5108 In: *Discrete Mathematics* 26.3 (1979), pp. 261–272.
- 5109 [186] Sergei Svitashhev, Joshua K Young, Christine Schwartz,
 5110 Huirong Gao, S Carl Falco, and A Mark Cigan. "Targeted
 5111 mutagenesis, precise gene editing, and site-specific gene
 5112 insertion in maize using Cas9 and guide RNA." In: *Plant*
 5113 *physiology* 169.2 (2015), pp. 931–945.
- 5114 [187] Yoshiyasu Takefuji and L Chen. "Parallel algorithms for
 5115 finding a near-maximum independent set of." In: *IEEE*
 5116 *Trans. Neural Networks* 1.3 (1990), p. 263.
- 5117 [188] Akito Taneda. "MODENA: a multi-objective RNA inverse
 5118 folding." In: *Advances and applications in bioinformatics and*
 5119 *chemistry: AACB* 4 (2011), p. 1.
- 5120 [189] Akito Taneda. "Multi-Objective Genetic Genetic for Pseu-
 5121 doknotted RNA Sequence Design." In: *Frontiers in Ge-*
 5122 *netics* 3 (2012), p. 36. ISSN: 1664-8021. doi: [10.3389/fgene.2012.00036](https://doi.org/10.3389/fgene.2012.00036). URL: <https://doi.org/10.3389/fgene.2012.00036>.
- 5125 [190] Akito Taneda. "Multi-objective optimization for RNA de-
 5126 sign with multiple target secondary structures." In: *BMC*
 5127 *bioinformatics* 16.1 (2015), pp. 1–20.

- 5128 [191] Michela Taufer, Abel Licon, Roberto Araiza, David Mireles, FHD Van Batenburg, Alexander P Gulyaev, and
 5129 Ming-Ying Leung. "PseudoBase++: an extension of Pseu-
 5130 doBase for easy searching, formatting and visualization
 5131 of pseudoknots." In: *Nucleic acids research* 37.suppl_1 (2009),
 5132 pp. D127–D135.
- 5134 [192] Siqi Tian and Rhiju Das. "RNA structure through multi-
 5135 dimensional chemical mapping." In: *Quarterly reviews of*
 5136 *biophysics* 49 (2016).
- 5137 [193] Pilar Tijerina, Sabine Mohr, and Rick Russell. "DMS foot-
 5138 printing of structured RNAs and RNA–protein complexes." In: *Nature protocols* 2.10 (2007), pp. 2608–2623.
- 5140 [194] Ignacio Tinoco Jr and Carlos Bustamante. "How RNA
 5141 folds." In: *Journal of molecular biology* 293.2 (1999), pp. 271–
 5142 281.
- 5143 [195] Ignacio Tinoco, Olke C. Uhlenbeck, and Mark D. Levine.
 5144 "Estimation of Secondary Structure in Ribonucleic Acids." In: *Nature* 230.5293 (Apr. 1971). Number: 5293 Publisher:
 5145 Nature Publishing Group, pp. 362–367. ISSN: 1476-4687.
 5146 DOI: [10.1038/230362a0](https://doi.org/10.1038/230362a0). URL: <https://www.nature.com/articles/230362a0> (visited on 04/14/2021).
- 5149 [196] Craig Tuerk and Larry Gold. "Systematic evolution of
 5150 ligands by exponential enrichment: RNA ligands to bac-
 5151 teriophage T4 DNA polymerase." In: *science* 249.4968
 5152 (1990), pp. 505–510.
- 5153 [197] Douglas H. Turner and David H. Mathews. "NNDB: the
 5154 nearest neighbor parameter database for predicting sta-
 5155 bility of nucleic acid secondary structure." In: *Nucleic
 5156 Acids Research* 38.suppl_1 (2009), pp. D280–D282.
- 5157 [198] Douglas H Turner and David H Mathews. "NNDB: the
 5158 nearest neighbor parameter database for predicting sta-
 5159 bility of nucleic acid secondary structure." In: *Nucleic
 5160 Acids Research* 38.suppl_1 (2010), pp. D280–D282.
- 5161 [199] Sinan Uğur Umu and Paul P Gardner. "A comprehensive
 5162 benchmark of RNA–RNA interaction prediction tools for
 5163 all domains of life." In: *Bioinformatics* 33.7 (2017), pp. 988–
 5164 996.

- 5165 [200] Jason G Underwood, Andrew V Uzilov, Sol Katzman,
 5166 Courtney S Onodera, Jacob E Mainzer, David H Math-
 5167 ewns, Todd M Lowe, Sofie R Salama, and David Haussler.
 5168 “FragSeq: transcriptome-wide RNA structure probing us-
 5169 ing high-throughput sequencing.” In: *Nature methods* 7.12
 5170 (2010), pp. 995–1001.
- 5171 [201] Gandhimohan M Viswanathan, EP Raposo, and MGE Da
 5172 Luz. “Lévy flights and superdiffusion in the context of
 5173 biological encounters and random searches.” In: *Physics*
 5174 of Life Reviews 5.3 (2008), pp. 133–150.
- 5175 [202] Alexey G Vitreschak, Dmitry A Rodionov, Andrey A
 5176 Mironov, and Mikhail S Gelfand. “Riboswitches: the old-
 5177 est mechanism for the regulation of gene expression?”
 5178 In: *Trends in Genetics* 20.1 (2004), pp. 44–50.
- 5179 [203] Manja Wachsmuth, Gesine Domin, Ronny Lorenz, Robert
 5180 Serfling, Sven Findeiß, Peter F Stadler, and Mario Mörl.
 5181 “Design criteria for synthetic riboswitches acting on tran-
 5182 scription.” In: *RNA biology* 12.2 (2015), pp. 221–231.
- 5183 [204] Haoyi Wang, Hui Yang, Chikdu S Shivalila, Meelad M
 5184 Dawlaty, Albert W Cheng, Feng Zhang, and Rudolf Jaenisch.
 5185 “One-step generation of mice carrying mutations in mul-
 5186 tiple genes by CRISPR/Cas-mediated genome engineer-
 5187 ing.” In: *cell* 153.4 (2013), pp. 910–918.
- 5188 [205] Shouhua Wang, Ting ting Su, Huanjun Tong, Weibin Shi,
 5189 Fei Ma, and Zhiwei Quan. “CircPVT1 promotes gallblad-
 5190 der cancer growth by sponging miR-339-3p and regulates
 5191 MCL-1 expression.” In: *Cell Death Discovery* 7.1 (2021),
 5192 pp. 1–10.
- 5193 [206] Zhong Wang, Mark Gerstein, and Michael Snyder. “RNA-
 5194 Seq: a revolutionary tool for transcriptomics.” In: *Nature*
 5195 reviews genetics
- 10.1 (2009), pp. 57–63.
- 5196 [207] Richard B Waring and R Wayne Davies. “Assessment
 5197 of a model for intron RNA secondary structure relevant
 5198 to RNA self-splicing—a review.” In: *Gene* 28.3 (1984),
 5199 pp. 277–291.
- 5200 [208] James D Watson and Francis HC Crick. “Molecular struc-
 5201 ture of nucleic acids: a structure for deoxyribose nucleic
 5202 acid.” In: *Nature* 171.4356 (1953), pp. 737–738.

- 5203 [209] Lina Weinbrand, Assaf Avihoo, and Danny Barash. "RNAf-
5204 binv: an interactive Java application for fragment-based
5205 design of RNA sequences." In: *Bioinformatics* 29.22 (2013),
5206 pp. 2938–2940.
- 5207 [210] Eric Westhof and Valérie Fritsch. "RNA folding: beyond
5208 Watson–Crick pairs." In: *Structure* 8.3 (2000), R55–R65.
5209 ISSN: 0969-2126. DOI: [https://doi.org/10.1016/S0969-2126\(00\)00112-X](https://doi.org/10.1016/S0969-2126(00)00112-X). URL: <https://www.sciencedirect.com/science/article/pii/S096921260000112X>.
- 5212 [211] Kay C Wiese, Andrew Hendriks, and Jagdeep Poonian.
5213 "Algorithms for RNA folding: a comparison of dynamic
5214 programming and parallel evolutionary algorithms." In:
5215 *2005 IEEE Congress on Evolutionary Computation*. Vol. 1.
5216 IEEE. 2005, pp. 475–483.
- 5217 [212] Kevin A Wilkinson, Edward J Merino, and Kevin M Weeks.
5218 "Selective 2'-hydroxyl acylation analyzed by primer ex-
5219 tension (SHAPE): quantitative RNA structure analysis
5220 at single nucleotide resolution." In: *Nature protocols* 1.3
5221 (2006), pp. 1610–1616.
- 5222 [213] Wade C Winkler and Ronald R Breaker. "Genetic con-
5223 trol by metabolite-binding riboswitches." In: *Chembiochem*
5224 4.10 (2003), pp. 1024–1032.
- 5225 [214] SA Woodson. "Recent insights on RNA folding mecha-
5226 nisms from catalytic RNA." In: *Cellular and Molecular Life*
5227 *Sciences CMLS* 57.5 (2000), pp. 796–808.
- 5228 [215] Xiufeng Yang, Kazuki Yoshizoe, Akito Taneda, and Koji
5229 Tsuda. "RNA inverse folding using Monte Carlo tree
5230 search." In: *BMC bioinformatics* 18.1 (2017), p. 468.
- 5231 [216] Hua-Ting Yao, Cedric Chauve, Mireille Regnier, and Yann
5232 Ponty. "Exponentially few RNA structures are designable."
5233 In: *Proceedings of the 10th ACM International Conference on*
5234 *Bioinformatics, Computational Biology and Health Informatics*.
5235 2019, pp. 289–298.
- 5236 [217] Hua-Ting Yao, Jérôme Waldspühl, Yann Ponty, and Se-
5237 bastian Will. "Taming Disruptive Base Pairs to Reconcile
5238 Positive and Negative Structural Design of RNA." In:
5239 *RECOMB 2021-25th international conference on research in*
5240 *computational molecular biology*. 2021.

- 5241 [218] Joseph N Zadeh, Conrad D Steenberg, Justin S Bois, Brian
 5242 R Wolfe, Marshall B Pierce, Asif R Khan, Robert M Dirks,
 5243 and Niles A Pierce. "NUPACK: Analysis and design of
 5244 nucleic acid systems." In: *Journal of computational chemistry*
 5245 32.1 (2011), pp. 170–173.
- 5246 [219] Joseph N Zadeh, Brian R Wolfe, and Niles A Pierce. "Nu-
 5247 cleic acid sequence design via efficient ensemble defect
 5248 optimization." In: *Journal of computational chemistry* 32.3
 5249 (2011), pp. 439–452.
- 5250 [220] Shay Zakov, Yoav Goldberg, Michael Elhadad, and Michal
 5251 Ziv-Ukelson. "Rich parameterization improves RNA struc-
 5252 ture prediction." In: *Journal of Computational Biology* 18.11
 5253 (2011), pp. 1525–1542.
- 5254 [221] Wenbing Zhang and Shi-Jie Chen. "RNA hairpin-folding
 5255 kinetics." In: *Proceedings of the National Academy of Sciences*
 5256 99.4 (2002), pp. 1931–1936.
- 5257 [222] Wenbing Zhang and Shi-Jie Chen. "Analyzing the biopoly-
 5258 mer folding rates and pathways using kinetic cluster
 5259 method." In: *The Journal of chemical physics* 119.16 (2003),
 5260 pp. 8716–8729.
- 5261 [223] Wenbing Zhang and Shi-Jie Chen. "Exploring the com-
 5262 plex folding kinetics of RNA hairpins: I. General fold-
 5263 ing kinetics analysis." In: *Biophysical journal* 90.3 (2006),
 5264 pp. 765–777.
- 5265 [224] Wenbing Zhang and Shi-Jie Chen. "Exploring the com-
 5266 plex folding kinetics of RNA hairpins: I. General fold-
 5267 ing kinetics analysis." In: *Biophysical Journal* 90.3 (2006),
 5268 pp. 765–777.
- 5269 [225] Qi Zhao, Zheng Zhao, Xiaoya Fan, Zhengwei Yuan, Qian
 5270 Mao, and Yudong Yao. "Review of machine learning
 5271 methods for RNA secondary structure prediction." In:
 5272 *PLoS computational biology* 17.8 (2021), e1009291.
- 5273 [226] Yu Zhu, ZhaoYang Xie, YiZhou Li, Min Zhu, and Yi-Ping
 5274 Phoebe Chen. "Research on folding diversity in statistical
 5275 learning methods for RNA secondary structure predic-
 5276 tion." In: *International Journal of Biological Sciences* 14.8
 5277 (2018), p. 872.
- 5278 [227] Michael Zuker and David Sankoff. "RNA secondary struc-
 5279 tures and their prediction." In: *Bulletin of mathematical
 5280 biology* 46.4 (1984), pp. 591–621.

- 5281 [228] Michael Zuker and Patrick Stiegler. "Optimal computer
5282 folding of large RNA sequences using thermodynamics
5283 and auxiliary information." In: *Nucleic acids research* 9.1
5284 (1981), pp. 133–148.
- 5285 [229] C. Zwieb. "Tmrdb (tmRNA database)." In: *Nucleic Acids
5286 Research* 28.1 (2000), pp. 169–170. doi: [10.1093/nar/28.1.169](https://doi.org/10.1093/nar/28.1.169).
5287 URL: <https://doi.org/10.1093/nar/28.1.169>.
- 5288 [230] C. Zwieb. "Tmrdb (tmRNA database)." In: *Nucleic Acids
5289 Research* 31.1 (2003), pp. 446–447. doi: [10.1093/nar/gkg019](https://doi.org/10.1093/nar/gkg019).
5290 URL: <https://doi.org/10.1093/nar/gkg019>.

[August 4, 2022 at 19:02 – 1.0]

5291 DECLARATION

5292 Put your declaration here.

5293 *Leipzig, June 2022*

5294

Nono Saha Cyrille Merleau

[August 4, 2022 at 19:02 – 1.0]

5296 COLOPHON

5297 This document was typeset using the typographical look-and-
 5298 feel `classicthesis` developed by André Miede and Ivo Pletikosić.
 5299 The style was inspired by Robert Bringhurst's seminal book on
 5300 typography "*The Elements of Typographic Style*". `classicthesis` is
 5301 available for both L^AT_EX and LyX:

5302 <https://bitbucket.org/amiede/classicthesis/>

5303 Happy users of `classicthesis` usually send a real postcard to
 5304 the author, a collection of postcards received so far is featured
 5305 here:

5306 <http://postcards.miede.de/>

5307 Thank you very much for your feedback and contribution.