

ALBERT LUDWIGS UNIVERSITY OF FREIBURG

MASTER THESIS

Multi-site RNA-RNA Interaction Prediction

Author:

Yogapriya Ayyanarmoorthy

Supervisor:

Dr. Martin Raden

Examiner:

*Prof. Dr. Rolf Backofen
(Bioinformatics, University of
Freiburg)*

*Prof. Dr. Sebastian A. Will
(AMIBio, Laboratoire
d'informatique de l'École
Polytechnique, IPP, France)*

A thesis submitted in fulfillment of the requirements
for the degree of Master of Science
in the Bioinformatics Group,
Department of Computer Science

Submitted on May 15, 2020

DECLARATION

I hereby declare, that I am the sole author and composer of my Thesis and that no other sources or learning aids, other than those listed, have been used. Furthermore, I declare that I have acknowledged the work of others by providing detailed references of said work. I hereby also declare, that my Thesis has not been prepared for another examination or assignment, either wholly or excerpts thereof.

Place , Date

Signature

I would like to thank Prof. Dr. Rolf Backofen so much for offering me this great opportunity to work on my Master Thesis in his group of Bioinformatics in Albert Ludwigs University of Freiburg.

I owe a lot to my supervisor Dr. Martin Raden, this thesis would not have been possible without his continuous guidance and support with both knowledge and ideas throughout the research. I want to thank him a lot for all his efforts, for providing me with thorough feedback and help throughout my work. Working with him was a great pleasure for me.

I am also grateful to my respectful parents for their great support along my stay, without whom I wouldn't have been able to stay abroad.

I want to thank my friends for being helpful and friendly, which contributed positively to my work.

Abstract

Ribonucleic acid (RNA) is an essential biological macromolecule in all biological cells. Computational prediction techniques may be used to determine how two RNA molecules will form intermolecular base pairing. A variety of biophysical and biochemical approaches are there to test the RNA-RNA interactions (RRIs). At very large computational expense, there are a range of algorithms in the literature that can anticipate a lot of these interactions.

The identification of non-coding (nc)RNA targets is largely regulated by two factors, namely the consistency of the duplex between the two interacting RNAs and the internal structure of both mRNA and ncRNA. Approaches may also be divided between various major categories based on whether they consider the inter- and intramolecular structure. One such approach is accessibility based interaction prediction model, which helps us in finding the single-site interaction of two RNAs. IntaRNA is a tool for rapid and precise prediction of interactions for accessibility based approach.

Within this thesis, we developed a model using an iterative scheme which predicts concurrent blocks of interactions within an accessibility based prediction model and provides us with the prediction of joint structure for the interacting RNAs and total energy. The proof for the energy of an respective Multi-site RRI are computed from the two energies is also shown. The respective extensions of the IntaRNA package will be included in the main package for external usage and further development. The comparison between various RNAs molecules have been tested and the output is been provided with polygon plots. Further to that, alternative RRIs prediction approaches and advance improvements and drawbacks are discussed within the thesis.

Zusammenfassung

Ribonukleinsäure (RNA) ist ein essentielles biologisches Makromolekül in allen biologischen Zellen. Computergestützte Vorhersagetechniken können verwendet werden, um zu bestimmen, wie zwei RNA-Moleküle eine intermolekulare Basenpaarung bilden. Es gibt verschiedene biophysikalische und biochemische Ansätze, um die RNA-RNA-Wechselwirkungen (RRIs) zu testen. Bei sehr hohem Rechenaufwand gibt es in der Literatur eine Reihe von Algorithmen, die viele dieser Wechselwirkungen antizipieren können.

Die Identifizierung von durch nicht-kodierende RNAs regulierten RNAs wird weitgehend durch zwei beeinflusst reguliert, nämlich die Konsistenz des Duplex zwischen den beiden interagierenden RNAs und die interne Struktur von mRNA und ncRNA. Ansätze können auch in verschiedene Hauptkategorien unterteilt werden, je nachdem, ob sie das inter- und intramolekulare Gerüst berücksichtigen. Ein solcher Ansatz ist das auf Accessibility basierende Interaktionsvorhersagemodell, das uns hilft, die Single-Site-Interaktion zweier RNAs zu finden. IntaRNA ist ein Werkzeug zur schnellen und präzisen Vorhersage von Interaktionen für mit Hilfe eines solchen Accessibility basierenden Ansatz.

In dieser Arbeit haben wir ein Modell entwickelt, das gleichzeitig auftretende Interaktionen mit Hilfe eines auf Accessibility basierenden Modells vorhersagt. Der Beweis für die Energie eines jeweiligen RRI mit mehreren Standorten, der aus den beiden Energien berechnet wird, wird ebenfalls gezeigt. Die jeweiligen Erweiterungen des IntaRNA-Pakets werden zur externen Verwendung und Weiterentwicklung in das Hauptpaket aufgenommen. Der Vergleich zwischen verschiedenen RNAs-Molekülen wurde getestet und mit Polygon-Plots visualisiert. Darüber hinaus diskutierten wir über den Vergleich zwischen den RRI-Vorhersageansätzen und über Verbesserungen und Nachteile.

Contents

| | |
|--|------------|
| Abstract | iii |
| Kurzfassung | iv |
| 1 Introduction | 1 |
| 1.1 Biological Background of RNA | 1 |
| 1.2 Formal background of RNA | 3 |
| 1.2.1 RNA Structure | 3 |
| 1.2.2 Nested secondary structure | 3 |
| 1.2.3 Nearest neighbor model and energy contributions | 4 |
| 1.2.4 Structure probabilities and McCaskill algorithm | 6 |
| 1.3 RNA-RNA Interaction | 7 |
| 1.3.1 Formal background of RNA-RNA interactions | 8 |
| 1.4 RNA-RNA Interaction Prediction Approaches | 11 |
| 1.4.1 Hybridization-only interaction prediction | 11 |
| 1.4.2 General RNA-RNA interaction prediction | 13 |
| 1.4.3 Concatenation-based RNA-RNA interaction prediction | 14 |
| 1.4.4 Accessibility-based interaction prediction | 15 |
| 1.4.5 Comparison of approaches for RRI prediction | 17 |
| 2 Multisite Accessibility Based | 19 |
| 2.1 RNAup - Exact Recursion for single site | 19 |
| 2.2 IntaRNA - Heuristic recursion for single site | 21 |
| 2.3 Iterative scheme for double-site RRI | 22 |
| 2.4 Generalization to multi-site RRI prediction | 24 |
| 3 Results & Discussion | 25 |
| 3.1 Setup | 25 |
| 3.2 OxyS – fhlA | 25 |
| 3.3 Spot42 – sthA | 28 |
| 3.4 GcvB – oppA | 32 |
| 3.5 DicF - ftsZ | 34 |
| 3.6 S-mRNA - EGS | 36 |
| 3.7 Details of studied RRI | 38 |
| 4 Summary | 39 |
| Bibliography | 40 |

| | |
|------------------------|-----------|
| Appendices | 43 |
| A RNA Sequences | 44 |

Chapter 1

Introduction

RNA molecules play important roles in various biological processes. Their regulation and function are mediated by interacting with other molecules, e.g by forming base pairs between two RNAs, called RNA-RNA interactions (RRI). Many RNAs interact via multiple synchronous, non-overlapping subinteractions (M-RRI), e.g. OxyS-fhlA. There are fast and reliable single interaction site (S-RRI) prediction tools like IntaRNA, which helps us in predicting mRNA target sites for given non-coding RNAs (ncRNAs), also they are capable of modelling all sites individually but not in a joint prediction. The simultaneous prediction of both intra- and inter-molecular base pairing allowing for multiple sites is computationally expensive. Some known approaches are IRIS, piRNA, NUPACK. Here we use a S-RRI prediction tool (namely IntaRNA) for the prediction of M-RRI.

1.1 Biological Background of RNA

In this thesis, I will focus on Ribonucleic acids (RNA). First of all, I would like to provide the basic biological background that is essential for the thesis. Ribonucleic acid, or RNA is one of the three major biological macromolecules that are important for all known forms of life (along with DNA (deoxyribonucleic acid) and proteins). The interaction between two RNAs plays a vital role in the basic cellular activities like transcription, RNA processing and translation. The process by which DNA is copied to RNA is called *transcription*, and that by which RNA is used to produce proteins is called *translation*. RNAs also play an important role in protein synthesis.

DNA is double stranded and RNA is a single-stranded molecule. Each strand of RNA is a sequence of four building blocks called *nucleotides*. Each nucleotide contains sugar, phosphate and nitrogen containing bases. The sugar and phosphate groups form the backbone of RNA strand and the bases bond to each other. The RNA molecules are represented as a sequence $S \in \{A, C, G, U\}^*$, where A (adenine), C (cytosine), G (guanine), U (uracil) are the bases of the nucleotide chain.

According to their potential for coding, RNA's are classified into two major categories i.e., coding RNAs and noncoding RNAs. Coding RNAs mostly refers to mRNA that encodes protein to act as different components including cell structures, signal transducers and enzymes. Non-coding RNAs act as cellular regulators with no protein encoding.

Complementary bases *C-G* and *A-U* form stable base pairs with each other using hydrogen bonds. These are called Watson-Crick pairs. Also important are the weaker *U-G* wobble pairs.

Together they are called *canonical base pairs*. In general, isolated base pairs are unstable. If interacting bases belong to the same molecule of RNA, they form *intra-molecular* structures and if they belong to different molecules of RNA, they form *inter-molecular* structures.

The prediction of RNA-RNA interaction is intended to predict these intermolecular structures between two RNA molecules, an extremely important step in understanding the role of ncRNAs. However, intramolecular and intermolecular structures are not mutually exclusive.

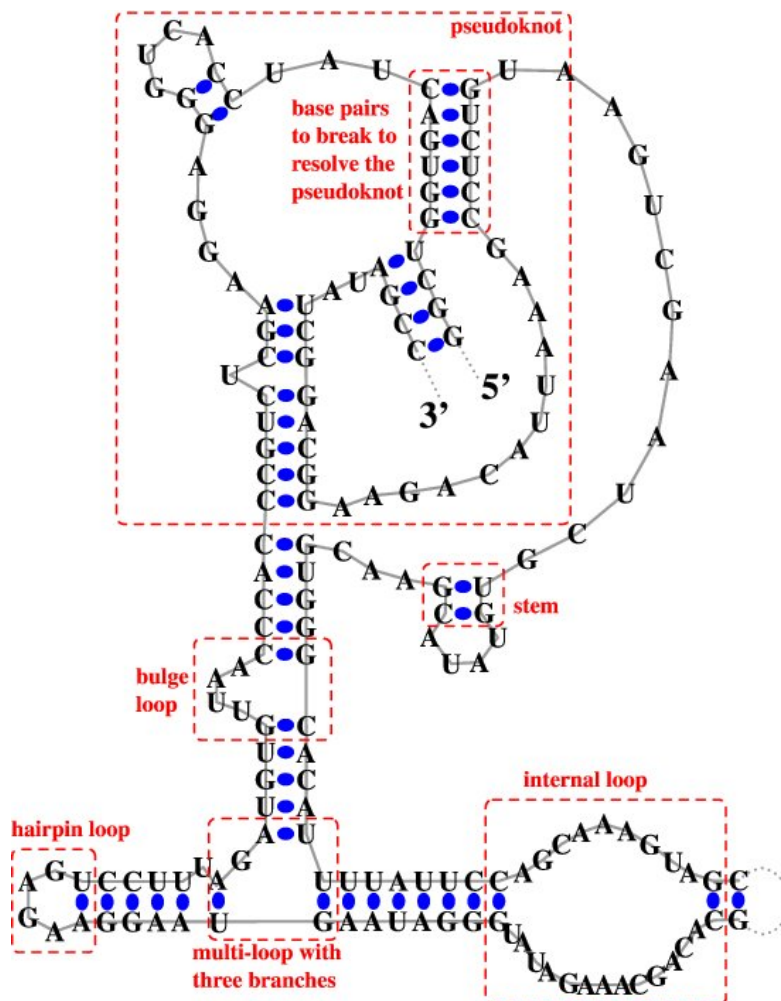


Figure 1.1: Schematic representation of the secondary structure (a set of base pairs) for the RNase P RNA molecule of *Methanococcus maripaludis* from the RNase P Database. Thick blue dots represents base pairs and red dashed boxes represent structural features such as stacking, bulges, hairpin, interior, multi loops and pseudoknot structure. This figure was taken from the RNAstrand webpage. (Andronescu et al., 2008)

Single stranded nucleic acid sequences contain many complementary regions that can form double helices when the molecule is folded back onto itself. The resulting pattern of double

helical stretches interspersed with loops is called the *secondary* structure of an RNA.

1.2 Formal background of RNA

Here in this section, I would like to bring up the formal definitions of ribonucleic acid.

1.2.1 RNA Structure

The RNA molecules are represented as a sequence $S \in \{A, C, G, U\}^*$. Formally, an RNA secondary structure P of S is a set of base pairs:

$$P \subseteq \{(i, j) | 1 \leq i < j \leq n, Si \text{ and } Sj \text{ are complementary}\},$$

where $n = |S|$ and for all $(i, j), (i', j') \in P$:

$$(i = i' \Leftrightarrow j = j') \text{ and } i \neq j'$$

To form a valid secondary structure, the base pairs must satisfy a number of limitations. Let the bases be numbered from 1 to n in a sequence. If the bases are complementary, a base pair may form between positions i and j , if $|j - i| \geq 4$, since there must usually be at least three unpaired bases in a hairpin loop. Let bases k and l form another allowed pair. The pair (k, l) is said to be compatible with the pair (i, j) if the two pairs can be present in a structure simultaneously. Pairs are compatible if they are non-overlapping (e.g. $i < j < k < l$) or if one is nested within the other (e.g. $i < k < l < j$). The Final case, where the pairs are interlocking or crossing (e.g. $i < k < j < l$) is called pseudo-knot. These pairs are assumed to be incompatible with most programs. An allowed secondary structure is a set of base pairs that are all compatible with each other.

They are different types of RNA secondary structure. i.e. nested and crossing structures. Crossing structures contain pseudo-knots, where two structure parts overlap. Nested structures doesn't have any crossing arcs.

1.2.2 Nested secondary structure

Nested secondary structures can be uniquely decomposed into so called loops or secondary structure elements. Depending on the number of enclosed base pairs and unpaired bases, different types of secondary structure elements are distinguished. These are hairpin loop, stacking, bulge loop, internal loop, multi loop.

Let S be a fixed sequence. Further, let P be an RNA structure for S .

- a base pair $(i, j) \in P$ is a *hairpin* loop if $\forall i < i' \leq j' < j : (i', j') \notin P$.
- a base pair $(i, j) \in P$ is a *stacking* if $(i + 1, j - 1) \in P$
- two base pairs $(i, j) \in P$ and $(i', j') \in P$ form an *internal* loop (i, j, i', j') if $i < i' < j' < j$; $(i' - i) + (j - j') > 2$; no base pair (k, l) between (i, j) and (i', j')

- An internal loop is called left (right, resp.) *bulge* if $j = j' + 1$ or $i' = i + 1$
- A *k-multiloop* consists of multiple base pairs, $(i_1, j_1) \dots (i_k, j_k) \in P$ with a closing base pair $(j_0, i_{k+1}) \in P$ with the property that $\forall 0 \leq l \leq k : (j_l < i_{l+1})$; $\forall 0 \leq l, l' \leq k$ there is no base pair $(i', j') \in P$ with $i' \in [j_l \dots i_{l+1}]$ and $j' \in [j_{l'} \dots i_{l'+1}]$.
- $(i_1, j_1) \dots (i_k, j_k)$ are called the *helices* of the multiloop.

1.2.3 Nearest neighbor model and energy contributions

DeVoe and Tinoco (1962) said vertical stacking of bases gives largest contribution to the stability of the RNA helix. The stacking of unpaired bases is less predictable and stable than the paired bases. Hence, the directly neighboured bases must be taken into account while estimating the energy contribution of a base pair, that results in the *Nearest Neighbor Model* (Borer et al., 1974).

The Nearest Neighbor Model enables the calculation of a free energy estimate for a given RNA secondary structure. The free energy can be taken as the amount of energy stored in a system. The system is more stable when the energy is lower. Hence, for the *most stable structure* of RNA, we go for *minimum free energy (MFE)*. The energy difference between the reference state to the system is measured. We have a reference system which we use to understand the stability of the system. The reference is an RNA structure with no base pairing (the open chain) ie., $E(\phi) = 0$. Hence, we need to check not only the hydrogen bonds but also the stacking stability. The Nearest Neighbor Model uses a loop-based structure decomposition. To avoid the duplication of stacking, only inner stacking are taken into account.

The terminal mismatch consists of the first unpaired bases immediately after the stacking. The identity of the terminal mismatch provides the energy of the loop. In Bulge or Internal loop also we have the same energy contribution. Energy contributions for external base pairs, which are not enclosed by any other base pairs, are referred to as *dangling end contributions*. The energy $E(P)$ Eq. 1.2.1 of a nested secondary structure P can be estimated by the sum of loop contributions (see Figure 1.2)

$$E(P) = \sum_{(i,j) \in P} \begin{cases} e^H(i, j) & : \text{if hairpin loop,} \\ e^{SBI}(i, j, k, l) & : \text{if stack/bulge/internal loop,} \\ e^M(i, j, x, x') & : \text{if Multi loop,} \end{cases} \quad (\text{Eq. 1.2.1})$$

Where e^H , e^{SBI} and e^M tells the context sensitive energy contributions of the loops. (k, l) represents the enclosed base pair of stack, bulge or internal and x represent the unpaired bases and x' represents the helices enclosed in the multi loop. There is an exponential number of possible multi loop composition. The energy for them can be estimated as below

$$e^M(i, j, x, x') = e_a^M + e_b^M x + e_c^M x'$$

where the pseudo energy parameter e_a^M scores the multi loop closing base pair (i, j) , e_b^M represents the penalty for x directly enclosed unpaired bases x and e_c^M scores x' enclosed helices .

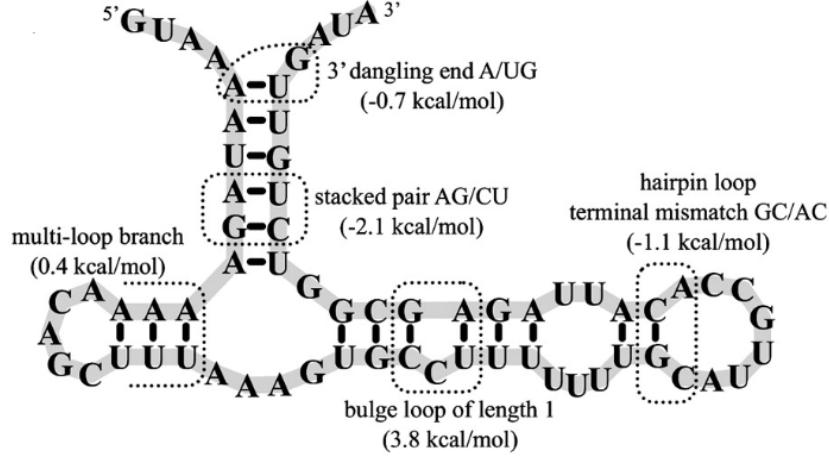


Figure 1.2: Energy contributions of loops. Figure is taken from (Andronescu et al., 2010)

Thus the nearest neighbor model gives the energy contributions for the loop types.

From the above energy model, We can define a recursive dynamic programming algorithm to compute the structure which minimizes the energy function, this is called minimum free energy (mfe) structure. This algorithm was introduced by Zuker and Stiegler (1981).

The basic substructures of the secondary structure of the RNA sequence (i.e., stack, hairpin, internal and multi loop) are independent of each other and the energy of the secondary structure is assumed to be the sum of the energies of the substructures. The algorithm is executed in two steps with a single RNA sequence as input. Firstly, the minimum free energy of the input RNA sequences is calculated, then traceback is used to recover the respective secondary structure with the base pairs. Thus given an RNA sequence S , Zuker's algorithm predicts the non-crossing, minimal energy structure P of S in $O(n^3)$ time and $O(n^2)$ space.

The most commonly used are the Turner parameters ((Mathews et al., 1999), (Mathews et al., 2004)) which can be found in the Nearest Neighbor Data Base NNDB. There are significant differences in the loop energies especially in the multi-loop scoring. Below is the parameter sets from the Vienna RNA package.

For multi-loop param (using turner 1999), by using e_{abc}^M as introduced above

$$A = 3.4, B = 0.4, C = 0$$

The stability of a 3 stem multi loop is calculated by the equation:

$$\Delta G = A + B + C$$

$$A + 2(B) \text{ \{loop degree of 2\}}$$

therefore, the value is 4.2

For multi-loop param (using turner 2004),

$$A = 9.3, B = -0.9, C = 0$$

The stability of a 3 stem multi loop is calculated by the equation:

$$\Delta G = A + B + C$$

$$A + 2(B) \text{ \{loop degree of 2\}}$$

therefore, the value is 7.5

1.2.4 Structure probabilities and McCaskill algorithm

Let's discuss about the structural information in terms of probabilities. According to the principal of maximum entropy (Jaynes, 1957) the best probability distribution for the calculation of the structure or base pair probability is the *Boltzmann Distribution*. These probabilities are calculated according to the Boltzmann weight. For RNA structures the unit of the energy value is $\frac{kcal}{mol}$ or $\frac{J}{mol}$. The RNA structure energy is been rescaled for Boltzmann weight computation. i.e., we replace Boltzmann constant k_B with the "mol-scaled" gas constant R to get the Boltzmann weight $w(P)$ of a structure P as:

$$w(P) = \exp\left(\frac{-E(P)}{RT}\right) \quad (\text{Eq. 1.2.2})$$

Where $E(P)$ represents the state energy, R represents the gas constant and T is the temperature.

The partition function Z can be calculated using the Boltzmann weights. Z is the sum of the Boltzmann weights of all states within \mathcal{P} , which is the set of all possible structures P that can be formed by S .

$$Z = \sum_{P \in \mathcal{P}} w(P) \quad (\text{Eq. 1.2.3})$$

Z is used for the calculation of structure and base pair probabilities. So in the total sum, the distribution does not change from a macroscopic point of view, therefore thermodynamic balance is reached.

The probability of an RNA structure P is given by

$$Pr[P|\mathcal{P}] = \frac{w(P)}{Z} \quad (\text{Eq. 1.2.4})$$

We can also calculate the probabilities of unpaired regions. Formally, we will identify the probability of the subsequences $i..j$ to be unpaired by $\mathcal{P}_{i,j}^u$. This probability depends on the whole ensemble of structures that can be formed by the RNA molecule of interest. Thus, it can be computed by

$$\mathcal{P}_{i,j}^u = \frac{Z_{i,j}^u}{Z}$$

where $Z_{i,j}^u$ is the partition function of all structures where the subsequence $i..j$ is unpaired. i.e.,

$$Z_{i,j}^u = \sum_{P \in \mathcal{P}_{i,j}^u} w(P) = Z(\mathcal{P}_{i,j}^u)$$

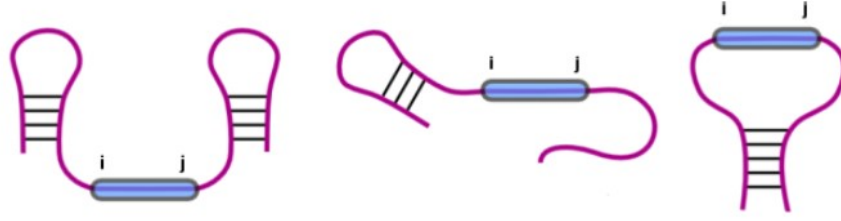


Figure 1.3: Exemplary structures that are unpaired in the subsequence $i..j$. The Figure was inspired by the lecture material of RNA bioinformatics lecture.

where $\mathcal{P}_{i,j}^u$ is the ensemble of all structures that are unpaired between i and j . i.e.,

$$\mathcal{P}_{i,j}^u = \{P \mid \nexists(k,l) \in P : i \leq k \leq j \text{ or } i \leq l \leq j\} \subseteq \mathcal{P}$$

The calculation of accessibility of single stranded regions is carried out using unpaired probability (Mückstein et al., 2006), hence it is very important.

Different probabilities can be calculated using McCaskill algorithm. The McCaskill algorithm (McCaskill, 1990) is used to calculate the partition function Z for a given sequence S , which can be used to compute probabilities. It enables efficient computing of the probabilities of the structure of the RNA as well as the probability that a certain base pair is formed. In addition, unpaired probabilities for subsequences can be calculated that reflect the accessibility of RNA parts for other interactions.

1.3 RNA-RNA Interaction

The interaction of RNA molecules is an essential factor for regulatory processes in all organisms. Computational prediction of RNA-RNA interactions (RRI) is a central methodology for the specific investigation of inter-molecular RNA interactions and regulatory effects of non-coding RNAs. RNA-RNA interactions are fast emerging as a major functional component in many newly discovered non-coding RNAs. They are important in many basic cellular activities including transcription, RNA processing, localization, and translation. Many RNA species function is guided by their structure, which is defined by intramolecular base pair formation. Small prokaryotic RNAs display evolutionary unstructured regions that control the expression of their target mRNAs by intermolecular base pairing (Wright et al., 2013). Hence, the prediction of both functional intramolecular and intermolecular structure of RNAs are important bioinformatics tasks.

Let's see about some simple RNA-RNA interactions. In *splicing*, small nuclear RNA's (snRNA) can recognize intronic regions of precursor messenger RNA(mRNA) which is the important step in identifying the RNA splicing products (Modrek and Lee, 2002). In *translation* transfer RNA(tRNA) interact with (mRNA) by "reading" the three letter code and define amino acid sequence (Selmer et al., 2006), (Ibba and Söll, 2000). In RNA modification, small nucleolar RNA(snoRNA) guide the modification of ribosomal RNA(rRNA) (Kiss, 2002). In microRNA (miRNA) targeting, the base pairing between miRNA and mRNA leads to degradation or translation inhibition of the mRNA (Bartel, 2004). For RNA function and regulation these examples

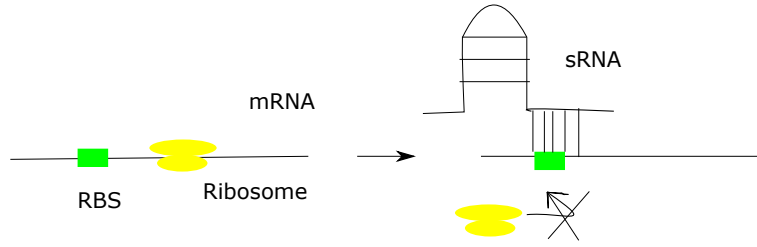


Figure 1.4: sRNA translation inhibition. sRNAs bind to ribosomal binding site (RBS) by complementary RNA-RNA interactions, which avoid ribosome binding.

show us the importance of the RNA-RNA interaction.

In order to allow highly accurate predictions, state-of-the-art methods not only take into account the stability (energy) of possible RNA-RNA interactions, but they also consider the accessibility of the interacting subsequences (Umu and Gardner, 2017), i.e., the intramolecular structure pattern.

Small regulatory RNAs (sRNAs) typically have a range of 50 to 500 nucleotides and are transcribed through intergenic bacterial genome regions. sRNAs work in trans and display complementarity with target mRNAs. Many sRNAs base pairs within the 5'-UTR of mRNA's target and blocks the ribosome binding site, therefore stops the initiation of translation. 1.4. This process is called Inhibition of translation initiation.

1.3.1 Formal background of RNA-RNA interactions

Here, we will see the formal background of RNA-RNA interactions.

In general RNA-RNA interaction prediction (RIP) problem, given two RNA sequences S^1 and S^2 (e.g., an antisense RNA and its target), the RIP problem asks one to predict their joint secondary structure. A joint secondary structure between S^1 and S^2 is a set of "pairings" where each nucleotide of S^1 and S^2 is paired with at most one other nucleotide, either from S^1 or S^2 (Alkan et al., 2006).

The RNA-RNA interaction is the combination of the set all of the base pairs in S^1 , the set of all base pairs in S^2 and the total intermolecular base pairs between two sequences. Formally, the RRI can be modelled as $\text{RRI} = \uplus \text{bp}(S^1) \cup \uplus \text{bp}(S^2) \cup \uplus (\text{Inter})$. Basically, the set of all base pairs of S^1 is P^1 and S^2 is P^2 , then *Inter* is I which denotes the set of all intermolecular base pairs.

$$\text{RRI} = P^1 \cup P^2 \cup I$$

Now, we further decompose I into the sequence of subsets of consecutive base pairs that form interaction blocks B which is depicted in the Figure 1.5, where $I = (B_1, \dots, B_x)$. A block B is the interaction block or interaction site.

Further, the interaction block or interaction site "B" can be represented as,

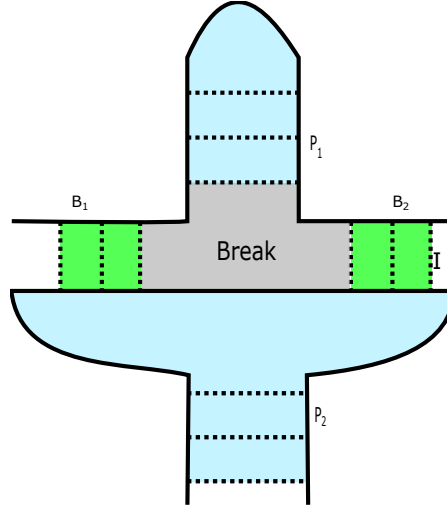


Figure 1.5: The RNA-RNA Interaction is the union of all base pairs in sequence 1 and sequence 2 are denoted as P^1 and P^2 (Blue colour) respectively. I (Green colour) denotes the union of all intermolecular base pairs. B_1, B_2 are the interaction blocks. Break (Grey colour) is the loop enclosed by two inter-molecular base pairs that also contains positions involved in intra-molecular base pairs

$$B = \{(i_1, i_2) \mid S_{i_1}^1 \text{ complementary to } S_{i_2}^2\} \subseteq [1, n_1] * [1, n_2]$$

Where for all $(i_1, i_2), (j_1, j_2)$ within a block B is

$$(i_1 < i_2) \iff (j_1 > j_2)$$

ie., They should be non-crossing. The block region $R(B)$ is $(i(B), j(B))$ ie., left and right most base pairs of B concerning S^1 .

$$i(B) = \arg \min_{i=(i_1, i_2) \in B} (i_1)$$



Figure 1.6: The left side figure shows the Positions are not paired within the loop. This problem starts with the pseudoknot which is shown in the right side figure where the same problem exists for the scoring of crossing structures.

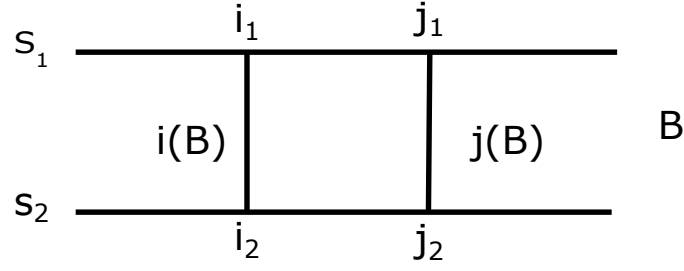


Figure 1.7: The block region $R(B)$ where the left and right most base pairs of B concerning S_1

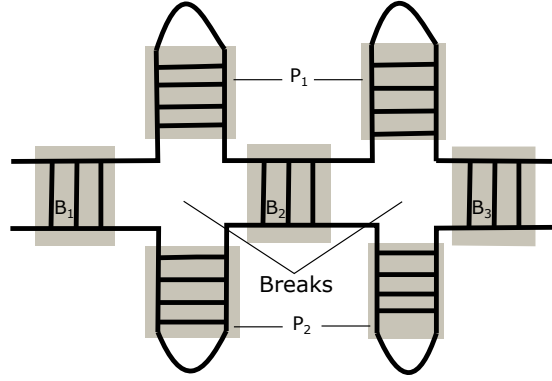


Figure 1.8: The interaction energy of RRI is the energy defined by the loops enclosed by all inter-molecular base pairs. $E(B_1)$ is the energy of block 1 and the $E(breaks)$ can be calculated from sum of all breaks.

$$j(B) = \arg \max_{i=(i_1, i_2) \in B} (i_1)$$

Furthermore, no intra molecular base pairs are allowed in block region $R(B)$ of P^1, P^2 . ie.,

$$\forall_{B \in I} : (\nexists_{(k,l) \in P^1} : i(B)_1 \leq k \leq j(B)_1 \vee i(B)_1 \leq l \leq j(B)_1) \wedge (\nexists_{(k',l') \in P^1} : i(B)_2 \leq k' \leq j(B)_2 \vee i(B)_2 \leq l' \leq j(B)_2)$$

where I is the union of all blocks (ie., all inter molecular base pairs) We compute the joint structure between S_1 and S_2 through minimizing their total free energy.

The Energy for the block $E(B)$ can be calculated as ,

$$E(B) = \sum_{\substack{i \in B \\ j = \arg \min(i') \\ i' < B \\ i'_1 > i_1}} E^{SBI}(i, j, k, l) \quad (\text{Eq. 1.3.1})$$

The $E(I)$ can be calculated as follows,

$$E(I) = E(\uplus B) + E_{init} \quad (\text{Eq. 1.3.2})$$

where, $E(\uplus B) = \sum_B E(B) + E(breaks)$ and E_{init} is fixed init score if $I \neq \phi$. E_{init} is 4.1 as per turner 99. Refer to section 1.2.3 .

In the Fig 1.5, light violet colour represents the intramolecular loop with the intermolecular base pairs paired. We will need to find out, how to score them. Here, without further knowledge or energy parameters, we score it via standard loop scores ignoring the intermolecular pairings. The problem is similar to pseudoknot scoring. These also contain the loops where the positions are not paired within the loop , see Figure 1.6.

$E(breaks)$ is defined by the sum over all individual breaks between blocks (Fig 1.8). For the $E(breaks)$ it depends on the prediction model which is a tricky part and that will be discussed with the next section along with the approaches idea. Now, we will summarise the formal definition of energy of RRI. By using the above RRI equation, we can write overall energy of RRI as

$$E(RRI) = E(P^1) + E(P^2) + E_{init} + \sum_{B_i \in I} E(B_i) + \sum_{B_i \in I} E_{break}(B_i, B_{i+1}, P^1, P^2) \quad (\text{Eq. 1.3.3})$$

1.4 RNA-RNA Interaction Prediction Approaches

There are several available methods, that can be classified according to their underlying prediction strategies, each implicating unique capabilities and restrictions often not transparent to the non-expert user.

Mostly for RNA-RNA interaction prediction methods are based on thermodynamic models and provide an efficient computation since Richard Bellman's principle of optimality (Raden et al., 2018) can be applied. RNA-RNA interaction prediction approaches are classified into hybrid-only interaction prediction, general interaction prediction, concatenation-based/cofolding interaction prediction, and accessibility-based interaction prediction.

In the following subsections we will see about the approaches used for predicting the RNA-RNA interactions.

1.4.1 Hybridization-only interaction prediction

In the hybrid-only interaction approach, the identification of RNA-RNA interaction doesn't consider intramolecular base pairs (fig 1.9) and they can be done with $O(nm)$ time and space complexity for two RNA sequences S^1 , S^2 of lengths n and m respectively (Tjaden et al., 2006).

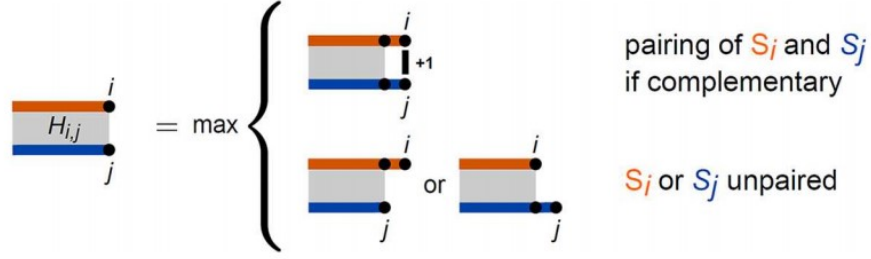


Figure 1.10: Recursion scheme to maximize intermolecular base pairs between two RNAs S_1 and S_2 represented in orange/blue, respectively

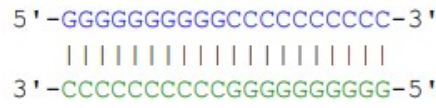


Figure 1.9: A full duplex structure where no intramolecular base pairs are assumed. The figure is taken from the paper (Wright et al., 2018)

A dynamic programming approach using a simplified energy model with two dimensional table H is filled via the prefix-based recursion Eq. 1.4.1.

$$H_{ij} = \max \begin{cases} H_{i-1, j-1} + 1 & : \text{if } S_i^1, \overleftarrow{S_j^2} \text{ are compl. base pair,} \\ H_{i-1, j} \\ H_{i, j-1}, \end{cases} \quad (\text{Eq. 1.4.1})$$

Where H_{ij} is the maximal number of intermolecular base pairs for the prefixes $S_1^1..i$ and $\overleftarrow{S_1^2}..j$ the reverse sequence of S^2 . The visual representation of the recursion scheme is given in Fig: 1.10. The above equation is the variant of the global sequence alignment approach by Needleman and Wunsch (1970) using scoring scheme i.e., base pair instead of match/mismatch scoring for S_i^1 , $\overleftarrow{S_j^2}$ no gap cost. Hence, when initialising $H_{i,0}/H_{0,j}$ with 0, the $H_{n,m}$ gives the maximal number of intermolecular base pairs and we can trace them back. As stated above, this approach has very low runtime, which are preserved when extended to energy minimization.

In order to compute the energy of an RRI using Eq. 1.4.1, no intra-molecular structure is considered, i.e. $P^1 = P^2 = \emptyset$.

Thus, eventually, only one block of inter-molecular base pairs is modelled i.e., $(I = B)$ and no break is present. They are implemented in tools like TargetRNA, RNAhybrid. The main advantage of these approaches is they are very fast and easy to calculate the significance of hits. Since, intramolecular base pairing is ignored they are used for the identification of short RNA's and otherwise they overestimate the length of target sites. These disadvantages can be overcome by concatenation and accessibility based approaches.

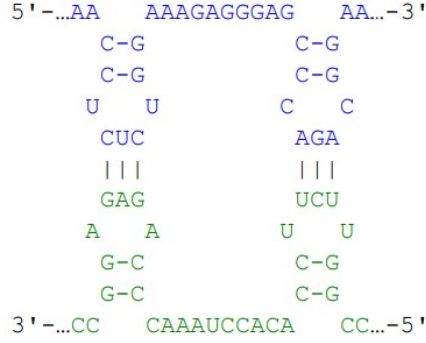


Figure 1.11: Double kissing hairpin interaction. The blue and green denotes the first and second sequence of RNA. Base pairs are denoted by dash. The picture is taken from (Wright et al., 2018)

1.4.2 General RNA–RNA interaction prediction

One of the most general approaches that is used for predicting the RRI of two intermolecular RNA molecules is IRIS (Pervouchine, 2004) method. This method is basically implemented by dynamic programming where it is the product of the sequence alignment and two MFOLD type secondary structure prediction algorithms. It can predict *general duplex structures*. This method is applied to some well known interactions such as OxyS with fhlA mRNA which basically forms a double kissing hairpin interactions as shown Fig. 1.11.

It shares most common features with pseudoknots, but is less computationally intensive. The input is made up of two sequences of RNA. Each sequence can form its own nested secondary structure and hybridize into the other molecule. The time and space complexities are $O(n^3m^3)$ and $O(n^2m^2)$, where n and m are the lengths of the sequences. The configuration of the OxyS-fhlA complex proposed in (Argaman and Altuvia, 2000) consists of four neighbouring stem loops, two in each of the molecules which connect, forming two stable kissing complexes. In this method, the main goal is the simultaneous optimization of intra- and inter-molecular base pairing.

IRIS also supports crossing of consecutive blocks treated by the last recursion case in the lower right of Fig. 1.12, which further complicates energy scoring of breaks. The energy contribution of general approach doesn't follow the interaction energy model instead they have pseudoknot energy. The energy associated with exterior pseudoknot can be given as (Xu and Chen, 2015)

$$G^{Pseudo} = \beta_1 + \beta_2 B^p + \beta_3 U^p \quad (\text{Eq. 1.4.2})$$

where β_1 represents penalty for introducing a pseudoknot, B^p is the number of base pairs that border the interior of the pseudoknot (i.e. number of paired positions) and U^p is the number of unpaired bases inside the pseudoknot 1.6(right). If a pseudoknot is inside a multiloop then they can be represented as β_1^m (by replacing the β_1) and if pseudoknot is inside another pseudoknot they can be represented as β_1^p (by replacing β_1).

As an approximation, one could use $E(PK-loop)$ with such pseudoknot energy terms based

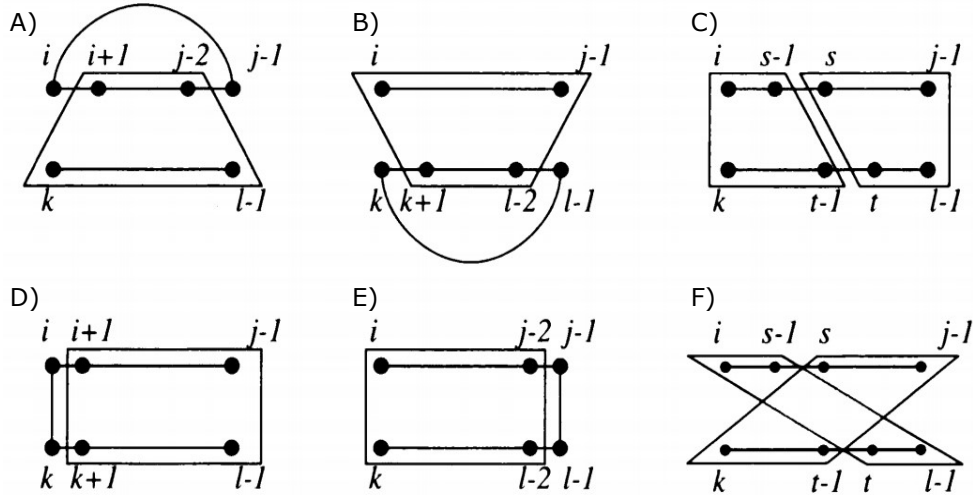


Figure 1.12: Depiction of the recursion $M_{j..l}^{i..k}$ which handles intramolecular (a,b) and inter-molecular (d,e) base pair extensions as well as a general decomposition (c) and crossing (f) case. Figure is taken from the paper Pervouchine (2004)

on G^{Pseudo} to score breaks. Note, to get an even more accurate overall energy scoring of an interaction, one would have to use pseudoknot energy terms also for such loops formed by intra-mol base pairs (refer 1.6 (left)). For simplicity, Eq. 1.3.3 uses only nested energy terms to assess intra-molecular energies. Thus, the exact energy computation of the general approach is not covered by the formalizations used within this thesis.

The time and space usage of IRIS are $O(n^6)$ and $O(n^4)$, respectively. The partition function version of RNA-RNA interaction prediction allows computation of probabilities of intermolecular interactions, which is used to access the stability. Due to its high complexity, several methods for reducing the requirements have been introduced. One such approach introduced by Chitsaz et al. (2009a) called *biRNA* is used to predict the multiple simultaneous binding site.

1.4.3 Concatenation-based RNA-RNA interaction prediction

Concatenation or co-folding approach is used for predicting the interacting base pairs of two RNA molecules based on intramolecular structure prediction methods. Here, two or more sequences are concatenated into a single sequence with special inter-spacing linker sequences. The final sequence is used within an adaptation of a standard structure prediction that takes care of the linker sequences. *mfold* was the first Concatenation-based prediction tool using the Nearest-Neighbor energy model, later implemented in *Mutli RNAfold* and *RNAcofold*.

RNA sequences are concatenated by a linker of length $l + 1$, where l is the minimal loop size, to ensure the concatenated sequence ends can form a base pair. We don't need any special energy treatment because the intra- and inter- molecular loops are treated equally. Hence the breaks are considered as multi loop and scored accordingly.

Concatenation-based approaches overcome the disadvantage of hybrid-only approach by in-



Figure 1.13: a) Pattern that can be predicated by Concatenation b)Kissing stem-loop and c) kissing hairpin interaction. Both (b) and (c) cannot be predicted as they form a crossing structure in the concatenated model. The blue and orange are the two different RNA's and the dotted green is the linker, black lines represents the base pairs. Figure inspired from paper (Raden et al., 2018)

corporating the competition of intra- and intermolecular base pairing. Still they cannot predict all interaction patterns because both intra- and intermolecular base pairs have to be nested. For example, interactions like kissing stem-loop or kissing hairpin-loop (as seen in fig 1.13) cannot be predicted by standard tools because they form a pseudoknot in the concatenation model.

NUPACK which is a pseudoknot prediction tool that solves the problem but with the higher runtime. They are based on dynamic programming approaches for specific classes of pseudoknot structures, but does not seem to be a significant drawback in terms of the accuracy of predictions for shorter sequences. (Dirks and Pierce, 2003).

1.4.4 Accessibility-based interaction prediction

To overcome the drawbacks of concatenation approaches, Accessibility approaches have been introduced. The main aim of this approach is to integrate ensemble properties of the single sequences that are necessary for the interaction. It can predict single site interaction patterns of two respective RNA subsequences. Tools like RNAup and IntaRNA implement this strategy. Here we have to dissolve intra molecular structure before the intermolecular interaction is formed. That is, in order to form a stable interaction of intermolecular base pairs, the intra molecular base pairs have to be opened/broken.

We can classify single-site RNA-RNA interactions based on the structural context of the respective subsequences, which are

- exterior - not enclosed by any base pair.
- hairpin loop - directly enclosed by a base pair.
- non-hairpin loop - subsequence enclosed by two base pairs forming a bulge, interior or multi-loop.

IntaRNA can predict single-site interactions within any structural context of the respective subsequences, but concatenation-based approaches can only predict exterior-exterior context combinations. Energy scoring differs from normal $E(RRI)$, since intra-molecular structure is

only considered implicitly via ensemble energies.

The term *ensemble* refers to the set of all secondary structures which can be formed through an RNA sequence. In an RNA sequence S , the accessibility energy of a region $[i, k]$ is determined by the energy difference (referred to as ED):

$$ED(i, k) = -(E^{all} - E_{i,k}^u) \quad (\text{Eq. 1.4.3})$$

Where E^{all} denotes the energy of the set of all possible secondary structures that can be generated by sequence S and $E_{i,k}^u$ denotes the energy of the ensemble of structures which have a single stranded area $[i, k]$.

The partition function is the total of all states \mathcal{P} over the Boltzmann factors. The energy of the ensemble E^{all} is

$$E^{all} = -RT \ln(Z) \quad (\text{Eq. 1.4.4})$$

The probability of unpaired regions can be used for calculating the accessibility penalty for an interval $[i, k]$, as shown below:

$$\begin{aligned} ED(i, k) &= -(E(\mathcal{P}) - E(\mathcal{P}_{i,k}^u)) \\ &= E(\mathcal{P}_{i,k}^u) - E(\mathcal{P}) \\ &= -RT \ln(Z_{i,k}^u) - (-RT \ln(Z)) \\ &= -RT \ln\left(\frac{Z_{i,k}^u}{Z}\right) \end{aligned}$$

Hence,

$$ED(i, k) = -RT \ln(\mathcal{P}_{i,k}^u) \quad (\text{Eq. 1.4.5})$$

$$ED_{i,k}^1 = -RT \cdot \log(\mathcal{P}_{i,k}^{u1}) \quad (\text{Eq. 1.4.6})$$

$$ED_{j,l}^2 = -RT \cdot \log(\mathcal{P}_{j,l}^{u2}) \quad (\text{Eq. 1.4.7})$$

Therefore, the alternative $E(RRI)$ formula for a single interaction block B is:

$$E(RRI = (B)) = E(B) + E_{init} + E^{all1} + ED_{i_1(B), j_1(B)}^1 + E^{all2} + ED_{i_2(B), j_2(B)}^2 \quad (\text{Eq. 1.4.8})$$

Since $ED = E^u - E^{all}$ substituting ED value in $ED + E^{all}$ gives,

$$E^u - E^{all} + E^{all} = E^u$$

$$E(RRI) = E(B) + E_{init} + E^{u1} + E^{u2} \quad (\text{Eq. 1.4.9})$$

Since E^{all} is constant for a given sequence, accessibility-based approaches only optimize $(E(B) + E_{init} + ED_B^1 + ED_B^2)$. To this end P^u values are precomputed.

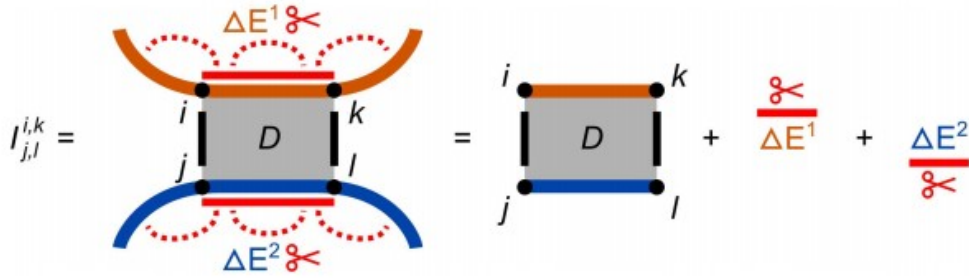


Figure 1.14: Depiction how accessibility-based approaches score an interaction of two RNAs $S1$ and $S2$ in orange and blue respectively. $\Delta E^1 + \Delta E^2$ are the energy needed to break the intramolecular base pairs and D is hybridization/duplex energy. Figure is taken from the paper (Raden et al., 2018)

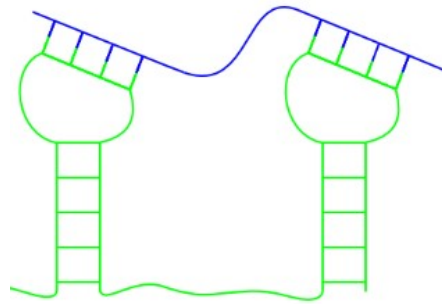


Figure 1.15: Double stem loop interaction cannot be handled by standard accessibility-based approaches as they have two binding sites (blocks) separated by intra- molecular structure. the figure is taken from COAT PhD summer school 2012

The energy of accessibility approach has no break , since the interaction I forms only one interaction block. Approaches like RNAup and IntaRNA use precalculated ED values for all possible interaction regions. They gives us how much energy is needed to free of intramolecular base pairs.

The main drawback of accessibility approach is, it can handle only one non-crossing block. These approaches cannot be modelled correctly for the double kissing hairpin interaction which has more than one crossing blocks of interaction 1.15

1.4.5 Comparison of approaches for RRI prediction

In this subsection, we will see the comparison between the approaches for some distinct interaction pattern. Below Table 1.1 gives an overview which interaction pattern can be predicted by which approaches.

Table 1.1: Comparison of RNA-RNA interaction prediction approaches for different interaction pattern

| Comparison of RRI approaches | | | | | | |
|------------------------------|--|--------------|---------------------------|-------------|------------------------|---------------|
| RRI Pattern | | | RRI prediction approaches | | | |
| Figures | RRI description | No.of blocks | Hybrid-only | General | Concatenation (nested) | Accessibility |
| 1.9 | Full duplex structure | 1 | yes | yes | yes | yes |
| 1.13 (a) | Nested joint structure without pseudoknots | 2 | no | yes | yes | no |
| 1.13 (b) | Stem loop interaction | 1 | no | yes | no | yes |
| 1.13 (c) | Kissing hairpin loop | 1 | no | yes | no | yes |
| 1.11 | Double kissing hairpin loop | 2 | no | yes | no | no |
| 1.5 | Kissing stem interaction | 2 | no | yes | no | no |
| 1.15 | Double kissing stem loop | 2 | no | yes | no | no |
| NA | Best Time complexity for each approach | NA | $O(nm)$ | $O(n^3m^3)$ | $O(n^3)$ | $O(n^2)$ |

The understanding of RNA structure and RNA-RNA interaction prediction approaches is important to ensure correct result interpretation and an knowing of their limitations are necessary to avoid wrong conclusions. Here we give a concise overview of the relevant theoretical history to the most general algorithmic approaches. We could say that the accessibility-based approach is the best approach for single site RNA-RNA interaction. To handle two or multi crossing blocks of interaction, we are introducing multisite accessibility based approach. The Multi-site RRI optimization is based on single-site IntaRNA predictions. Hence, we are going for the multisite accessibility based approach in the next chapter.

Chapter 2

Multisite Accessibility Based

In this chapter, we will introduce an accessibility-based approach that can be used for multisite RNA-RNA interaction prediction. In simple words we could say, it is Multi-site RRI optimization based on single-site IntaRNA predictions. Accessibility-based RNA-RNA interaction prediction methods are typically modelling a single block of consecutive inter-molecular base pairs. Thus, interaction pattern that consists of multiple concurrently formed blocks can not be predicted. Within this thesis, we are developing and testing possibilities to efficiently predict concurrent blocks of interaction within an accessibility-based prediction model. The approach will be based on IntaRNA, which is one of the state-of-the-art programs for RNA-RNA interaction prediction.

IntaRNA, developed by (Busch et al., 2008) and (Raden et al., 2018) at the bioinformatics group at Freiburg University, is a general and fast approach to the prediction of RNA-RNA interactions incorporating both the accessibility of interacting sites as well as the existence of a user-definable seed interaction. IntaRNA uses energy minimisation to find the best possible interaction.

Many RNAs interact via multiple synchronous, non-overlapping subinteractions (M-RRI), e.g. OxyS-fhlA. The simultaneous prediction of both intra- and inter-molecular base pairing allowing for multiple sites is computationally expensive. Current approaches include IRIS (Perovuchine, 2004), NUPACK (Dirks et al., 2007), piRNA (Chitsaz et al., 2009b), etc. There are fast and reliable single interaction site (S-RRI) prediction tools like RNAup and IntaRNA, that often show the additional sites within their suboptimal list, ie. are capable of modelling all sites individually but not in a joint prediction. To overcome this, we use the iterative method in this thesis for finding the interaction between multiple blocks. Beforehand, details of the S-RRI accessibility-based tools RNAup and IntaRNA are introduced.

2.1 RNAup - Exact Recursion for single site

In the following, I will first introduce the RNAup-like exact recursions (Mückstein et al., 2006) and then give an overview of IntaRNA heuristic version. The total energy score of the interaction is measured as the sum of the free hybridization energy and the free energy required to make the interaction sites available.

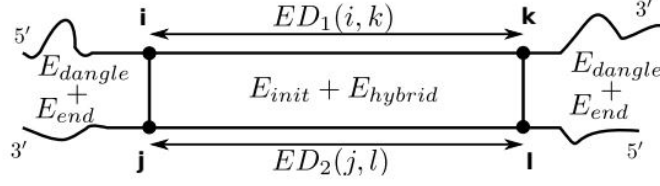


Figure 2.1: The energy contribution of IntaRNA. The image is taken from (Gelhausen, 2018)

Thus, Scoring an interaction in IntaRNA is dependent on two energy contributions:

- **Hybridization energy** : energy value E_{hybrid} from intermolecular base pairings in the form of stackings, bulges or internal loops, i.e., energy is typically a negative value.
- **Accessibility energy** : An amount of energy ED needed to single-strand the interacting region, i.e. not include intramolecular pairings, i.e., energy is a positive value.

The energy of an ensemble of structures is calculated using a partition function (McCaskill, 1990). Similarly, we get $ED(i, k)$ by calculating the partition function, $Z_{i,k}^u$ covering the ensemble of all structures which can be formed by a sequence S , with a single stranded region $[i, k]$. As introduced in the section 1.4.4. Therefore,

$$ED(i, k) = -RT \log(\mathcal{P}_{i,k}^u)$$

Mückstein et al. (2006) gives more detailed information on the same. The hybridization energy is measured using the Nearest Neighbor Energy Model. This represents the minimum free hybridization energy of two subsequences, where a base pair is generated by the left and right most positions of both subsequences. For sub-sequences $S_i^1 \dots S_k^1$ and $S_j^2 \dots S_l^2$, where S^1 is ordered from 5' to 3' and S^2 in the reverse order:

$$H(i, j, k, l) = \min_{\substack{B \\ i(B)=(i,j) \\ j(B)=(k,l)}} (E(B))$$

The hybridization energy is calculated with a Zuker-like recursion.

$$H(i, j, k, l) = \begin{cases} E_{init} & : \text{if } (S_i^1, S_j^2) \text{ can pair and } i = k, j = l, \\ \min_{r,s} \{e^{SBI}(i, j, r, s) + H(r, s, k, l)\} & : \text{if } (S_i^1, S_j^2) \text{ and } (S_k^1, S_l^2) \text{ can pair } i < k \text{ and } j < l, \\ \infty & : \text{otherwise,} \end{cases} \quad (\text{Eq. 2.1.1})$$

Here e^{SBI} is the energy contribution of stack, bulge and internal loop. The traceback helps us to find the base pairs of optimal interaction with energy $H(i, j, k, l)$. Both the accessibility and

hybridisation energy forms the extended hybridisation energy which is the specific hybridisation between $S_i^1 \dots S_k^1$ and $S_j^2 \dots S_l^2$ given by,

$$C(i, j, k, l) = \begin{cases} H(i, j, k, l) + ED_1(i, k) + ED_2(j, l) \\ \quad : \text{if } (S_i^1, S_j^2) \text{ and } (S_k^1, S_l^2) \text{ can pair } i \neq k \text{ and } j \neq l, \\ \infty \\ \quad : \text{otherwise,} \end{cases} \quad (\text{Eq. 2.1.2})$$

We get the time and space complexity of $O(n^2m^2)$ by limiting the loop size, which is still very high. When we limit the interaction length to l , it has a complexity of $O(nml^2)$ time and $O(nml^2)$ space. The interaction with the minimum estimated free energy (mfe) is probably the most stable structure and thus the structure fulfills the RNA molecule function. Therefore, we are interested in

$$mfe = \min_{i,j,k,l} C(i, j, k, l)$$

2.2 IntaRNA - Heuristic recursion for single site

The exact recursions of RNAup are not suitable for larger genome wide studies due to its high time and space complexity ie., $O(n^2m^2)$ where n represents the length of query and m is the length of target sequence. In order to overcome the time and space complexity problem, IntaRNA introduced the heuristic recursion. This recursion is based on sparsification technique. Hence, we consider only one right end of interactions with left end (i, j) which is single and locally optimal, instead of all the possible interaction. This will help us to reduce the space and time complexity to $O(nm)$, as introduced for IntaRNA version 1 & 2. The heuristic version is defined as:

$$C(i, j) = \min \begin{cases} E_{init} + ED_1(i, i) + ED_2(j, j) \\ \quad : \text{in the case of new interaction} \\ \min_{p,q} \{ e^{SBI}(i, j, p, q) + C(p, q) - ED_1(p, K(p, q)) - ED_2(q, L(p, q)) \\ \quad \quad + ED_1(i, K(p, q)) + ED_2(j, L(p, q)) \} \\ \quad : \text{if } (S_i^1, S_j^2) \text{ can pair ,} \\ \infty \\ \quad : \text{otherwise,} \end{cases} \quad (\text{Eq. 2.2.1})$$

Here $K(p, q), L(p, q)$ are newly introduced matrices which provide the right end of the best interaction with left end (p, q) . Since ED values are not additive, we have to subtract the old ED values before we add the new ED value.

IntaRNA also enforces a seed region as the required constraint for the interaction of two RNAs. Seed regions are a feature found in many RNA-RNA interactions. The seed region is an interaction region of almost complete complementarity. For animal microRNAs the seed region were first discovered (Bentwich, 2005), (Brennecke et al., 2005). Then, Tjaden et al. (2006) discovered it for many bacterial sRNAs.

RNAup does not use any seed condition. Within IntaRNA predictions, at least one seed is needed at the interaction site. The seed length is normally assumed to be between six and eight nucleotides. To speed up the genome-wide, methods such as RIssearch2 ((Alkan et al., 2017)) or RIBlast ((Fukunaga and Hamada, 2017)) utilize suffix-array-dependent screens to classify seed regions that are eventually expanded in both directions utilizing a predictive method focused on usability. In IntaRNA the length of the seed can be set by the user (Busch et al., 2008).

That is, given a set of putative seed interactions \mathcal{B}_{seed} . We introduce additional matrix C^{seed} to form a seed interaction. The seed matrix C^{seed} can be filled in using the following recursion:

$$C^{seed}(i, j) = \min \begin{cases} \min_{p, q} \{ e^{SBI}(i, j, p, q) + C^{seed}(p, q) - ED_1(p, K^{seed}(p, q)) - ED_2(q, L^{seed}(p, q)) \\ \quad + ED_1(i, K^{seed}(p, q)) + ED_2(j, L^{seed}(p, q)) \} \\ \quad : \text{if } (S_i^1, S_j^2) \text{ can pair ,} \\ \min_{p, q} \{ seed(i, j, p, q) + C(p, q) - ED_1(p, K(p, q)) - ED_2(q, L(p, q)) \\ \quad + ED_1(i, K(p, q)) + ED_2(j, L(p, q)) \} \\ \quad : \text{if } (S_i^1, S_j^2) \text{ can pair and are left end of the seed,} \\ \infty \\ \quad : \text{otherwise,} \end{cases} \quad (\text{Eq. 2.2.2})$$

The mfe of IntaRNA is

$$mfe = \min_{i, j} C^{seed}(i, j)$$

2.3 Iterative scheme for double-site RRI

In the following, we will introduce a new approach that uses a Single-RRI prediction tool (namely IntaRNA) for the prediction of Multi-RRI. For simplicity, the approach is first introduced for two sites B_1 and B_2 . To this end, an iterative scheme is to be applied, which is described in the following steps.

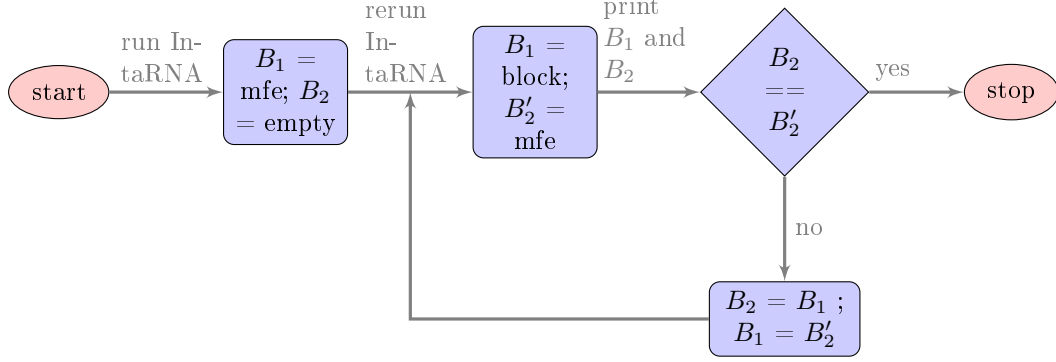
- Step 1: Firstly, we have to run IntaRNA and store minimum free energy and boundaries of respective B_1 .
- Step 2: Then we get the 'blocking' constraint from step 1 to rerun IntaRNA and predict the conditional minimum free energy and site B_2 . Here we block B_1 (Constrained IntaRNA) both for intra- and inter molecular base pairing and we get B_2 as minimum free energy. Then, the energy of a respective M-RRI can be computed from the two energies, ie., the energy of the conditional call can be added.

$$E(B_1 \wedge B_2) = E(B_1) + E(B_2|B_1) .$$

- Step 3: Since prediction of B_2 is conditional, the existence of B_2 can have effects B_1 . Thus, one starts to iterate the procedure from (2) but swaps the conditional site and check for convergence: is the site from two-steps before retained ($B_2 == B_2'$)? If yes:

convergence and stop iteration by printing the B_1, B_2 . If no: repeat constraint prediction until convergence by swapping .

Below is the flowchart representation for the same procedure.



Below is the proof of why the energy of the conditional call can be just added. The joint site has energy

$$\begin{aligned}
 E(B_1 \wedge B_2) &= E_{hyb}(B_1 \wedge B_2) + ED(B_1 \wedge B_2) \\
 &= E_{hyb}(B_1) + E_{hyb}(B_2) - RT\log(\mathcal{P}^u(B_1 \wedge B_2))
 \end{aligned} \tag{Eq. 2.3.1}$$

The first block B_1 is scored by

$$\begin{aligned}
 E(B_1) &= E_{hyb}(B_1) + ED(B_1) \\
 &= E_{hyb}(B_1) - RT\log(\mathcal{P}^u(B_1))
 \end{aligned} \tag{Eq. 2.3.2}$$

and the conditional prediction of B_2 by

$$\begin{aligned}
 E(B_2|B_1) &= E_{hyb}(B_2|B_1) + ED(B_2|B_1) \\
 &= E_{hyb}(B_2|B_1) - RT\log(\mathcal{P}^u(B_2|B_1))
 \end{aligned} \tag{Eq. 2.3.3}$$

Now, we add right end side values of Eqn Eq. 2.3.3 + Eq. 2.3.2, we get,

$$E_{hyb}(B_1) + E_{hyb}(B_2|B_1) - RT\log(\mathcal{P}^u(B_1)) - RT\log(\mathcal{P}^u(B_2|B_1)) \tag{Eq. 2.3.4}$$

As we know $\log(A) + \log(B) = \log(A * B)$, we apply this condition for log values in Eq. 2.3.4

$$\begin{aligned}
 &-RT\log(\mathcal{P}^u(B_1)) - RT\log(\mathcal{P}^u(B_2|B_1)) \\
 &= -RT\log(\mathcal{P}^u(B_1) * \mathcal{P}^u(B_2|B_1))
 \end{aligned}$$

Since $P(A \wedge B) = P(A) * P(B|A)$ and $E_{hyb}(B_2|B_1)$ is independent of B_1 , We get,

$$E_{hyb}(B_1) + E_{hyb}(B_2) - RT\log(\mathcal{P}^u(B_1 \wedge B_2)) \quad (\text{Eq. 2.3.5})$$

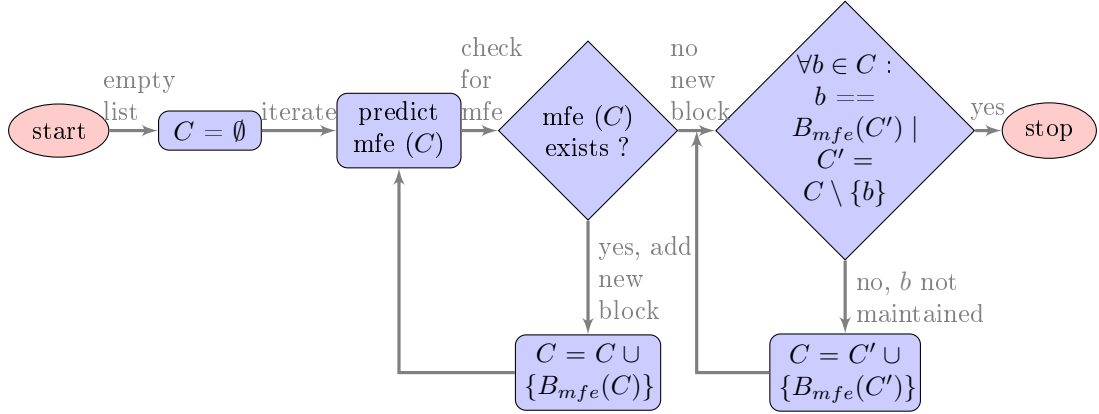
Now, we see the equations Eq. 2.3.5 and Eq. 2.3.1 are equal.

2.4 Generalization to multi-site RRI prediction

The generalized multi-site RRI prediction is similar to double-site RRI prediction. Here the process takes two steps (i) Iterative Accumulation, (ii) Iterative Refinement.

- Step 1: Iterative Accumulation : Here, We iterate the empty list of constrained prediction and predict the minimum free energy (mfe). Every iteration returns us a list of blocks C . Check for mfe (C) exists? then, we add a new block $C = C \cup \{B_{mfe}(C)\}$. Empty list of constraints are iterated until we don't get a new block. $C' = \{C\}$. In the end, we get a list of constraints.
- Step 2: Iterative Refinement :For every new block (ie., $\forall b \in C : b = B_{mfe}(C')$ where $C' = C \setminus \{b\}$), we need to check if it is preserved or not. If preserved, for all C , then we stop the iteration. If not, then we need to alter the set of blocks ie., $C = C' \cup \{B_{mfe}(C')\}$.

Below is the flowchart representation for the same procedure.



The mfe is

$$mfe^*(C) = \sum_{b \in C} (E_{hyb}(b) + E_{init}) + ED(C)$$

where $ED(C)$ is $-RT\log(P_u(C))$

Chapter 3

Results & Discussion

In order to evaluate the multi site interaction, I am comparing the results of the new approach for same RNA interactions with single site interaction tool IntaRNA and results from the literature. Details of the reported RRI are given at the end of the section in 3.2

3.1 Setup

I have used the IntaRNA-3.1.3-windows-64bit version for my thesis. The following parameter has been set. The `-outMode=C` a flexible interface to generate RNA-RNA interaction output in CSV format (using `;` as separator). The argument `-n 1` or `-outNumber=1` can be used to generate up to N interactions for each query-target pair. IntaRNA provides the possibility to constrain the accessibility computation using the `-qAccConstr` and `-tAccConstr` parameters. In this I have used "b" blocked to indicate the positions are occupied by some other interaction (implies single-strandedness). It is possible to restrict the overall length an interaction is allowed to have. This can be done independently for the query and target sequence using `-qIntLenMax` and `-tIntLenMax`, respectively. We can alter indexing (independently for query and target) using the `-qIdxPos0` and `-tIdxPos0` parameters, respectively. Here, the overall energy E has times of E_{init} .

With the above setup, we were able to test multi-site RNA-RNA interactions. The used sequences are listed in Appendix Table.

3.2 OxyS – fhfA

The small RNA OxyS binds to a short sequence inside the fhfA mRNA coding region. This is one of the classic examples of multi-site RRI where OxyS forms a stable kissing hairpin complex with fhfA.

Details of the interaction are taken from the paper (Argaman and Altuvia, 2000). The pairing mechanism between the two RNAs is dramatically influenced by their structure. For this purpose, a full comprehension of the pairing process involves thorough knowledge of the individual RNA structures.

When OxyS binds *fhlA* mRNA it forms kissing complex structure and results in a healthy anti-sense-target structure. The secondary structure of the 5' end region of *fhlA* mRNA was predicted to include two stem loop structures. The findings of the study of the structure confirm the presence of the two structure.

When we run with IntaRNA first, we observed that the OxyS (98 - 104) interacts with *fhlA* (-9 - -15), ie., 104:-15 & 98:-9 is the block that is predicted with the energy -4.37. Then, we run them with our tool, the total energy for both blocks is -7.99. Here the sequence length for OxyS starts from 1 to 109 and for *fhlA* it is from -53 to 60. We can clearly see from the Fig 3.2 where the predicted interaction (orange) and the original (blue) interaction from the paper interacts almost at the same place.

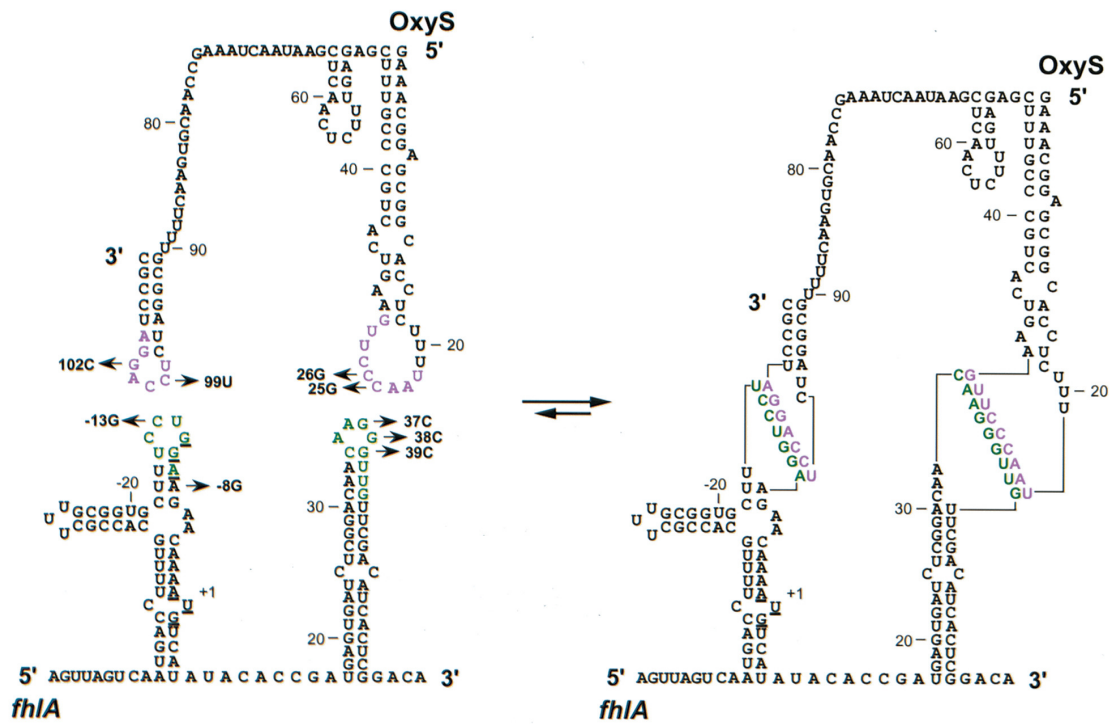


Figure 3.1: Interaction of OxyS - *fhlA*. The numbering of *fhlA* begins with the initiation codon (AUG). The counting of OxyS begins at the transcription start site. Figure is taken from (Argaman and Altuvia, 2000).

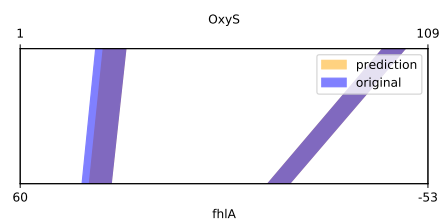


Figure 3.2: OxyS:fhfA interaction that is predicted using double-site tool. We can see that predicted interaction (orange colour) is highly similar to the original interaction (blue colour).

3.3 Spot42 – sthA

The next interesting example is Spot42 interaction with sthA. The RNA Spot42 plays a large role in the suppression of catabolites in Escherichia coli (E.coli) by the direct suppression of genes involved in primary and secondary metabolism.

This example is taken from (Beisel and Storz, 2011). Spot42 is interacting with its targets via three conserved accessible regions (I – III) refer to left side Fig 3.3 .

The mutation in region III influenced the repression of fusion of sthA the most (refer right side Fig 3.4). Mutating sites I and III were reported to have the highest impact while Region II mutation showed only slight effects. It has been shown that mainly sites I and III are important for the interaction with the target mRNA encoded by the sthA gene. This suggests multiple interaction possibilities of Spot42 with sthA, as discussed in Mann et al. (2017).

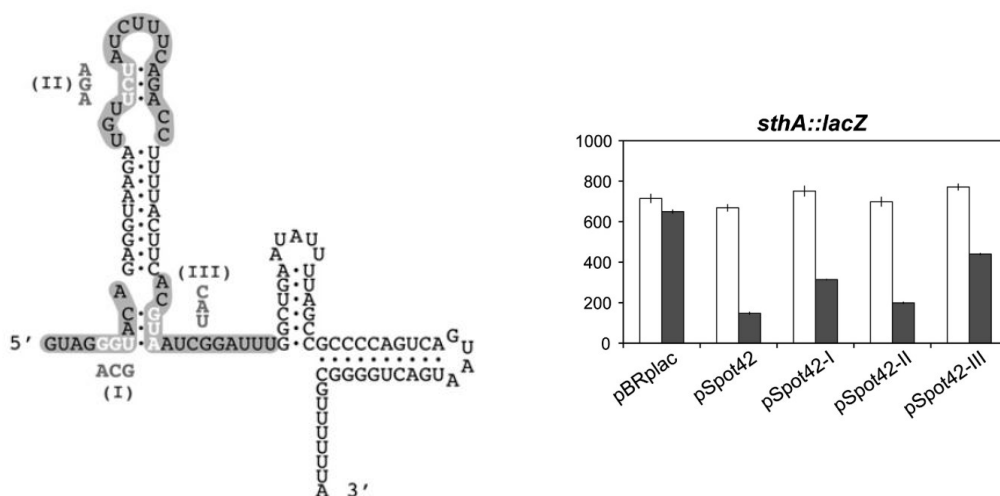


Figure 3.3: The left side of the figure shows the secondary structure of Spot 42. Mutated regions are in grey colour. Right side of the figure shows the mutational analysis of base-pairing interactions of Spot42 – sthA. Spot 42 base pairs directly with target genes via three separate regions is shown. Figure is taken from the paper (Beisel and Storz, 2011). The black bar shows the translation product measure of sthA gene and the white one is a control molecule to compare for the different experiments (translational reference).

When we run with IntaRNA first, we observed that the Spot42 (34 - 55) interacts with sthA (15 - 40), ie.,55:15&34:40 is the block that is predicted which refers to region III and the mfe is -7.85. Then, we run them with our tool here when Spot42 interacts with sthA, we observed that the total energy value is -26.38. The prediction provides a model for the concurrent interaction of region I and III. If we need to have the same ED values, then we can use the -energyNoDangles option. Here for Spot42, we couldn't find the original interaction, hence we are comparing with the bar chart from the figure 2 from paper (Beisel and Storz, 2011). From the right side figure 3.3, we can clearly see, mutating sites I and III showed the highest effect while region II mutation showed only small effects. This is because, I and III are close in structure though they are far away within the sequence.

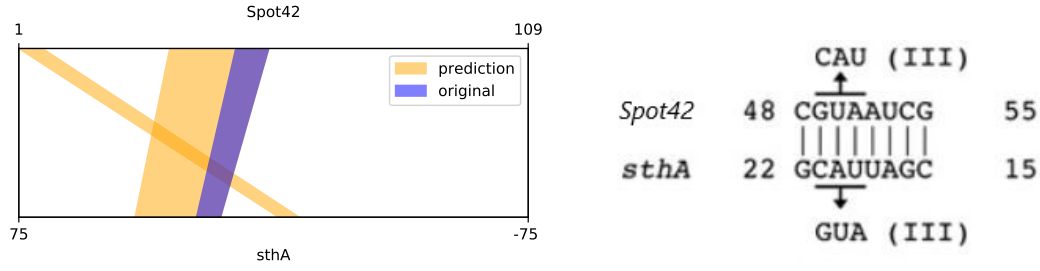


Figure 3.4: The left side of the figures shows the Spot42:sthA interaction that is predicted using double-site tool. We can see that predicted interaction (orange colour) of two blocks and only one original interaction block (blue colour) of original is shown as the other interaction is not given in the literature. The right side of the figure shows the base-pairing interactions with Spot 42 predicted by the folding algorithm NUPACK. Mutation in Region III is shown. Figure is taken from the paper (Beisel and Storz, 2011)

As our model shows only two interacting sites, we tried manually generating the third-site for that we run the IntaRNA with seedbp=6 while blocking the two first interaction sites (I and III). As a result, we get the interaction with site (II) (ie., Spot42 (20 - 25) interacts with sthA (-19 - -14)) with the mfe of -2.02. Thus the total interaction energy for all the three blocks is -15.07. Fig 3.5 shows that all three sites can contribute to the interaction as observed within the experiment of right side of Fig 3.3. As the interaction sites I and II are close to the start codon, they block the RBS which are called Inhibition of translation initiation. Refer to the Fig 1.4 for the sRNA translation inhibition.

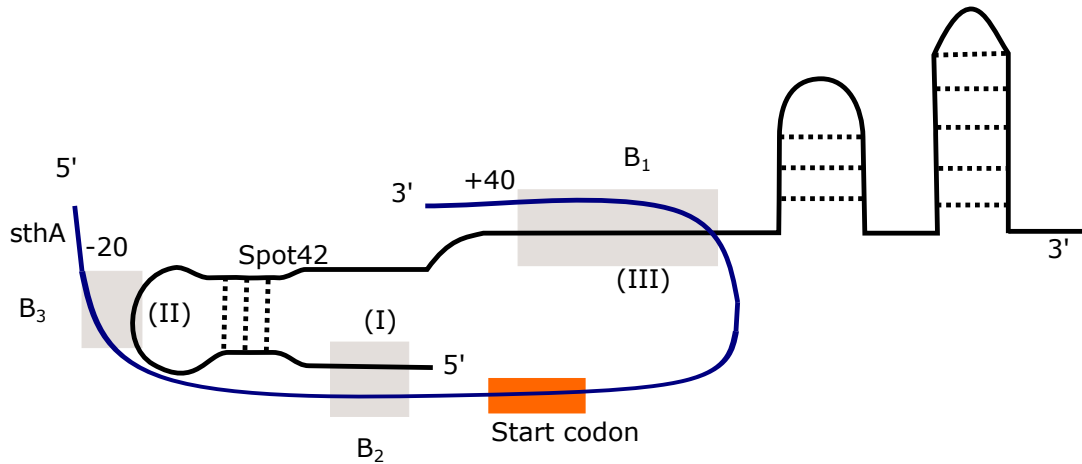


Figure 3.5: Spot42 interacts with sthA has three interacting sites (grey colour). As, our model shows only two interacting sites, we manually generated the third-site by blocking the interacting sites I and III and got the third interacting site II. As the start codon (orange colour) is very close to the interacting sites I and II, they stop the RBS. The Figure shows the all the three interaction sites of Spot42 (black colour) with sthA (blue colour).

Similarly, we manually generated the third site for the few other target mRNAs (gltA, srlA, nanC, xylF). The mutation in region I affected the repression of gltA, xylF and nanC. From the Figure 3.6 (except the top right one) we can also say that mutation site I showed the highest effect compared to the interaction site II and III for gltA, xylF and nanC. The mutation in region II affected the repression of srlA and from the Figure 3.6 (top right) we can also say that mutation site II showed the highest effect. The mutation in region III affected the repression of sthA.

Now, we tried manually getting the third-site. For gltA, we got the interaction site II (ie., Spot42 (25 - 30) interacts with gltA(80 - 85)) with the mfe of -0.47. Thus the total interaction energy for all three blocks of gltA is -12.21. For srlA, we got the interaction site I (ie., Spot42 (1 - 7) interacts with srlA(-166 - -160)) with the mfe of -3.94. Thus the total interaction energy for all three blocks of srlA is -17.25. For nanC, we got the interaction site II (ie., Spot42 (26 - 31) interacts with nanC(-160 - -155)) with the mfe of -1.09. Thus the total interaction energy for all three blocks of nanC is -21.31.

xylF results were quite different from other target mRNAs. When we were running this RNA in our tool, we got the block 1 interacting in mutation site I and II 3.7 (bottom right). The block 2 in mutation site III. Though xylF is interacting in both the mutation site I and II, it has its highest effect in mutation site I compared to site II. After that when we tried running for the third site manually using IntaRNA tool, we didn't get the third site. Thus the total interaction energy for two blocks of xylF is -14.15. These findings indicate that the three single-stranded regions of Spot 42 are involved in strong base-pair interactions with gltA, nanC, xylF, srlA and sthA mRNAs. Below table 3.1 gives you the details of each Spot42 target RNAs, which blocks covers which the mutation sites.

| Target mRNAs | mutation site I | mutation site II | mutation site III | Order of blocks |
|--------------|-----------------|------------------|-------------------|-------------------|
| sthA | B_2 | B_3 | B_1 | $B_3 < B_1 < B_2$ |
| gltA | B_1 | B_3 | B_2 | $B_1 < B_3 < B_2$ |
| srlA | B_3 | B_1 | B_2 | $B_2 < B_3 < B_1$ |
| nanC | B_1 | B_3 | B_2 | $B_1 < B_3 < B_2$ |
| xylF | B_1 | B_1 | B_2 | $B_1 < B_2$ |

Table 3.1: Interactions between the single stranded regions of Spot 42 and some target mRNAs. This table shows us for each Spot42 target RNA, which block covers which mutation site, along with the order of blocks in target.

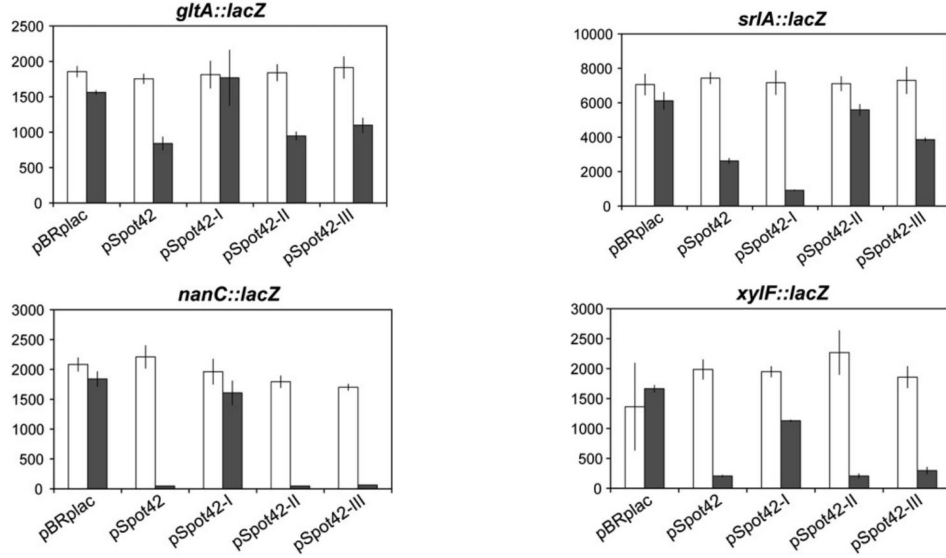


Figure 3.6: The figure shows the mutational analysis of base-pairing interactions of Spot42 with its target mRNAs. The black bar shows the translation product measure of *sthA* gene and the white one is a control molecule to compare for the different experiments (translational reference). Figure is taken from the paper (Beisel and Storz, 2011) .

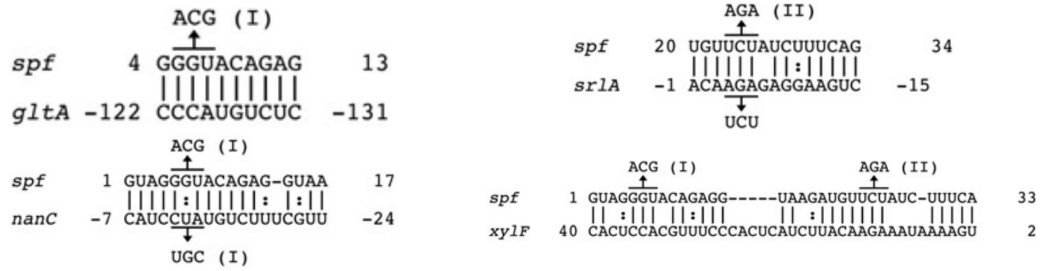


Figure 3.7: The figure shows the base-pairing interactions with Spot 42 predicted by the folding algorithm NUPACK. It shows the mutation region sites for the target mRNAs. Here *spf* is referred to as Spot42. Figure is taken from the paper (Beisel and Storz, 2011)

3.4 GcvB – oppA

The *gcvB* gene encodes two small, nontranslated RNAs that regulate *oppA*. The structure of the GcvB-*oppA* complex consists of two intermolecular helices that precede and follow the putative terminator. This is an example where there are four concurrent blocks predicted by the IRIS tool (Pervouchine, 2004).

oppA is the periplasmic-binding protein portion of the OppABCDF oligopeptide transport system. The functional consequence of deleting the *gcvB* gene is a derepression of *oppA*. The mechanism of GcvB regulation of *oppA* is likely to be on a translational level (Urbanowski et al., 2000). One of the major role of *oppA* is the transport of dietary peptides.

Study of the GcvB sequence identified a complementarity area near the ribosome-binding sites of *oppA* mRNAs. The findings from (Pulvermacher et al., 2008) indicate that various regions of GcvB have specific functions in the control of *oppA* mRNA. The Shine-Dalgarno sequence in the GcvB-*oppA* complex is obstructed (Pervouchine, 2004) and the complex structure is located in the upstream region. This is very much in accordance with the assumption that the *oppA* control seems to be at the translational stage.

When we run with IntaRNA first, we observed that the *oppA* (-9 - 14) interacts with GcvB (67 - 89), ie., 14:67&-9:89 is the block that is predicted with an mfe of -14.57. Then, we run them with our tool as, we predict the double block interaction, we have predicted the two sites of the interaction out of four. The total interaction energy for two sides of block is 26.38 kcal/mol. Also, in the prediction we get a crossing structure, which is in contrast to the model by IRIS , Fig : 3.8. Second block not in accordance with IRIS prediction as they form a pseudoknot. So, here the block predicted is not as same as the original.

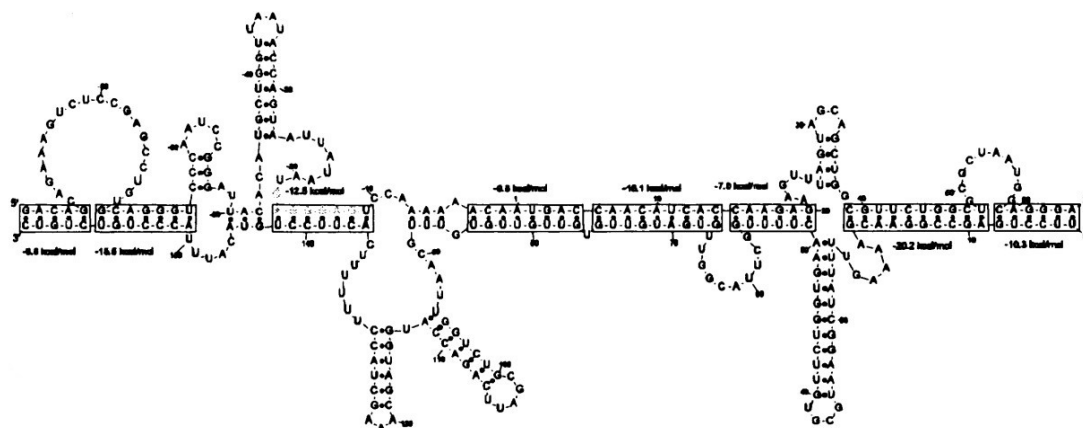


Figure 3.8: Interaction of GcvB – *oppA* . Figure is taken from the paper (Pervouchine, 2004)

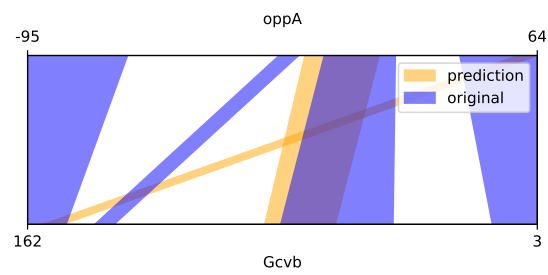


Figure 3.9: GcvB:oppA interaction that is predicted using double-site tool. We can see that out of four original interaction blocks(blue colour), we predicted (orange colour) the interaction of two blocks.

3.5 DicF - ftsZ

This is an example of a crossing interaction. We note that the DicF-ftsZ complex admits an area of complementarity, which gives rise to a generalized pseudo-knot.

The function of DicF was predicted on the basis of the complementarity of DicF RNA with the ftsZ mRNA binding region. DicF RNA is a 53-nucleotide gene formed in certain E.coli mutants. Here, DicF RNA is an antisense regulator of ftsZ translation.

This example is taken from the paper (Pervouchine, 2004). They say that DicF RNA has substantial complementarity with ftsZ mRNA in the area surrounding the Shine Dalgarno sequence, which is compatible with the finding that DicF controls ftsZ through interaction with the ribosome binding. DicF-ftsZ complex admits the region of complementarity which leads to generalized pseudoknot (refer to Fig 3.10).

When we run with IntaRNA first, we observed that the Dicf (35 - 52) interacts with ftsZ (53 - 75) is the block that is predicted (52:55&35:73) with mfe of -6.89. Then, we run them with our tool. For this example we used the parameter file and set the tIntLenMax=20 (for restricting the overall length an interaction).The total interaction energy for two sides of block that has been predicted by the tool is -13.86. In the prediction model, the cross structure is not formed, which says that the prediction model is different from the original one's.

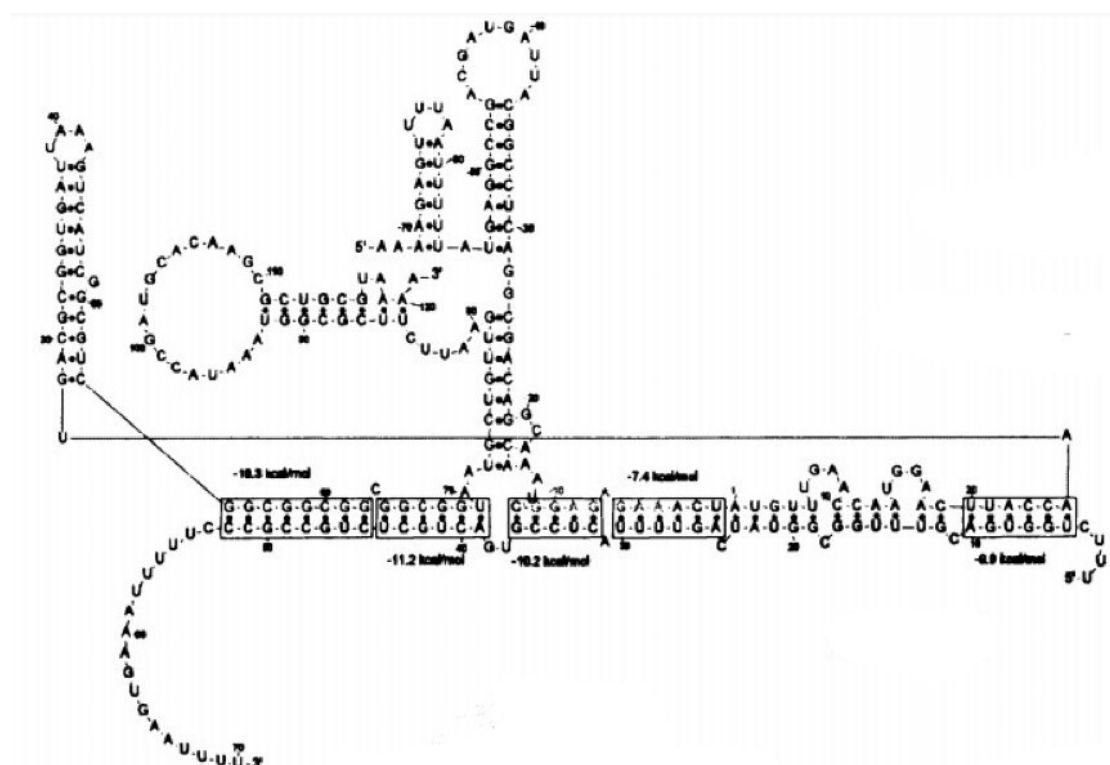


Figure 3.10: Interaction of DicF - ftsZ which has a generalized pseudoknot structure that has been predicted by IRIS. Figure is taken from the paper (Pervouchine, 2004).

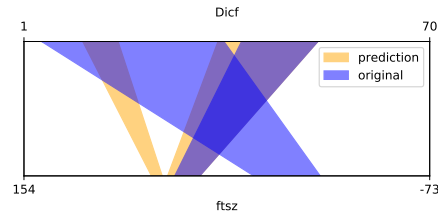


Figure 3.11: DicF:ftsZ interaction that is predicted using double-site tool. We can see that the original interaction blocks (blue colour) form a crossing structure, whereas the predicted (orange colour) blocks doesn't.

3.6 S-mRNA - EGS

In all species ribonuclease P (RNase P) was found. External guide sequences (EGSs) are RNA molecules consisting of a sequence complementary to mRNA targeting and recruiting intracellular ribonuclease P (RNase P), a tRNA processing enzyme, for target mRNA specific degradation. EGS RNAs derived from natural tRNA sequences can be good in blocking gene expression in bacteria.

This example is taken from the paper (Zhang et al., 2013). It is possible that an improvement in the RNase P cleavage rate could be attributed to additional tertiary interactions that theoretically stabilize the mRNA-EGS complex. Variant C386 was chosen for this analysis because the EGS RNAs derived from this version are among the most efficient EGS's. EGS S-C386 was built by connecting the EGS domain of C386 to targeting sequences complementary to the S mRNA. The EGS, S-SER, originating from the normal tRNA^{Ser} series, was also built. If this is the case, the binding affinity of the EGS variant (i.e. S-C386) to the target S RNA sequence might be greater than that of the EGS (i.e. S-SER) derived from the normal tRNA sequence. We couldn't reproduce the structures from literature for any of the sequence pairs.

S-C386-C and S-SER-C were derived from S-C386 and S-SER, respectively, and incorporated simple substitutions (5'-UUC-3' → AAG) at the three closely conserved locations in the T-loop of these EGSs. Nucleotides in these three positions are highly conserved among tRNA molecules and are essential for the folding and recognition of tRNA molecules by RNase P, so mutations in these positions are involved in the EGS process. S-C386-C and S-SER-C had the same anti-sense pattern to the target S RNA series as S-C386 and S-SER Fig: 3.12 and had identical binding affinities to S38 as S-C386 and S-SER, respectively. S-C386-C and S-SER-C can also be used as anti-sense regulation of such EGSs.

Due to its very short sequence Fig: 3.12, we took the 50nt long on the left side and right side of UCUUCAUCCUGCUGCUAUGCCUCAUCUUC of S-mRNA. When we run with IntaRNA first, we observed that the mRNA (-1 - 7) interacts with EGS S-SER & SER-C (46 - 53) i.e., -1:53&7:46 is the block that is predicted for both mRNA:EGS S-SER and mRNA:EGS S-SER-C with mfe of -8.41. The total energy for mRNA:EGS S-C386 is -14.14, mRNA:EGS S-C386C is -13.24, mRNA:EGS S-SER and mRNA:EGS S-SER-C is -16.0. Here, the start of the S-mRNA is complementary to the end of the S-SER, which led to the crossing pattern.

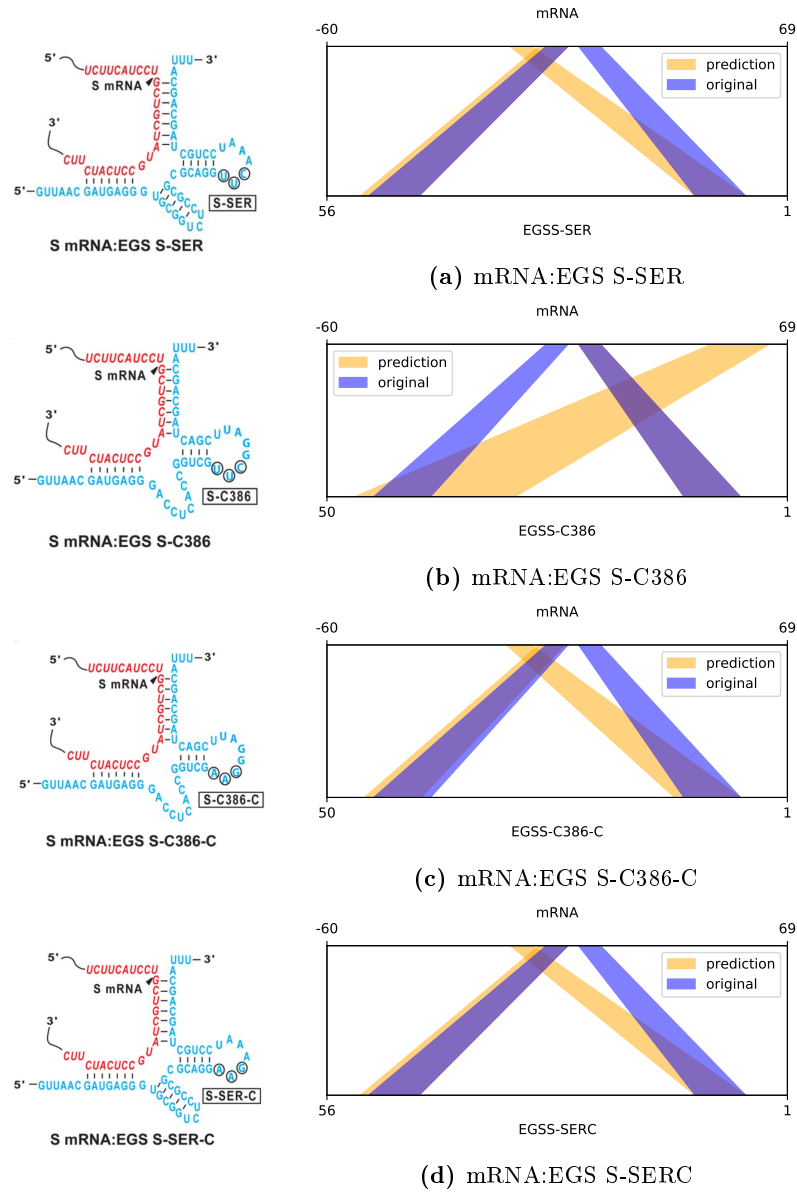


Figure 3.12: An EGS resembling the structure of a tRNA. The left side of the figure is taken from the paper (Zhang et al., 2013). The site of cleavage by RNase P is marked with an arrowhead, S mRNA (red color) and the EGS (blue color). The right side is the predicted structure from the tool. Original interaction blocks (blue colour) and the predicted blocks (orange colour). Just the exact sequence of the S mRNA around the targeting region is displayed (in red) and the EGS sequence is shown in blue. The RNase P cleavage site is labelled with an arrowhead. The sequences of S-SER and S-SER-C relative to T-stem and loop and the variable area of the tRNA molecule were derived from tRNASer, while those of S-C386 and S-C386-C were derived from the EGS variant C386. The index position is taken from the site of cleavage.

3.7 Details of studied RRI

Below is the table for the comparison between the different interacting RNAs that are used for the study purpose within this thesis. The table provides us with the original interaction from the paper and the prediction's from the tool. The table also clearly tells us that $E(B1 + B2)$ the multi-site energy is almost double the single site energy $E(B1)$. 1st Index gives us with the index positions and the length determines the sequence length.

| query:target | original | prediction | $E(B1)$ | $E(B1+B2)$ | 1st index | length |
|-------------------|---|------------------------------|---------|------------|-----------|---------|
| OxyS:fhlA | 104:-15&98:-9 30:34&22:42 | 104:-15&98:-9 30:24&24:40 | -4.37 | -7.99 | 1:-53 | 109:113 |
| Spot42:sthA | 55:15&48:22 7:NA&5:NA | 55:15&34:40 7:-8&1:-2 | -7.85 | -13.05 | 1:-75 | 109:150 |
| oppA:GcvB | -95:163&-64:151 -17:142&-11:136 -3:84&19:49 39:18&64:1 | 14:67&-9:89 63:152&57:158 | -14.57 | -26.38 | -95:3 | 159:160 |
| DicF:ftsZ | 52:55&39:69 36:-12&5:25 | 52:55&35:73 12:82&18:76 | -6.89 | -13.86 | 1:-73 | 70:227 |
| mRNA:EGSS-SER | 16:7&10:13 7:46&1:52 | -1:53&7:46 -9:13&-3:7 | -8.41 | -16.0 | -60:1 | 129:156 |
| mRNA:EGSS-SER-C | 16:7&10:13 7:46&1:52 | -1:53&7:46 -9:13&-3:7 | -8.41 | -16.0 | -60:1 | 129:156 |
| mRNA:EGS S-C386 | 1:46&7:40 10:13&16:7 | 47:48&63:31 10:13&16:7 | -8.15 | -14.14 | -60:1 | 129:150 |
| mRNA:EGS S-C386-C | 1:46&7:40 10:13&16:7 | -1:47&6:41 -10:14&-3:7 | -8.67 | -13.24 | -60:1 | 129:150 |

Table 3.2: Collections of multisite RNA interaction. The query and target are the RNAs interacting. ":" is used to differentiate between the query and target , "&" is used for the differentiate between start and end of the blocks.

Chapter 4

Summary

The motivation of the thesis is to efficiently predict concurrent blocks of interaction within an accessibility-based prediction model. Here, we use a Single-site RNA-RNA interaction prediction tool (namely IntaRNA) for the prediction of Multi-site RNA-RNA interaction. Due to the time constraint, we have implemented the double-site interaction here by using an iterative scheme. In the iterative scheme we block the interaction site and fixing the conditional site while running the IntaRNA tool. Until we get the convergence, we iterate. The theoretical approach of multi-site interaction scheme has also been discussed. Future work on the multi-site interaction can be implemented based on the theoretical approach that has been discussed above (refer to section 2.4).

The proof and implementation of double-site interaction model has been developed within this thesis. The theory proof of energy of an respective Multi-site RNA-RNA interaction can be computed from the two energies was developed and shown in the section 2.3. Then the idea has been implemented by using python based on IntaRNA calls. The double-site interaction study says that interaction energy for two blocks B1 and B2 together often double the minimum free energy. The detailed study of such examples are given in a table 3.2 for the easy comparison.

We collected some sample Multi-site RNA-RNA interaction examples from literature. The examples with the crossing interaction (pseudoknot), with two intermolecular helices, kissing interaction, etc., are used for the study purpose. We used the python polygon plots for plotting the blocks that are predicted. Then, we compared the outcome with extracted data from literatures. Each example shows its own structure, but most of them were very well close to the original one's that is taken from the literature. The brief explanation of the each study has been given in the chapter 3.

Bibliography

- Can Alkan, Emre Karakoc, Joseph H Nadeau, S Cenk Sahinalp, and Kaizhong Zhang. RNA–RNA interaction prediction and antisense RNA target search. *Journal of Computational Biology*, 13(2):267–282, 2006.
- Ferhat Alkan, Anne Wenzel, Oana Palasca, Peter Kerpedjiev, Anders Frost Rudebeck, Peter F Stadler, Ivo L Hofacker, and Jan Gorodkin. Risearch2: suffix array-based large-scale prediction of rna–rna interactions and sirna off-targets. *Nucleic acids research*, 45(8):e60–e60, 2017.
- Mirela Andronescu, Vera Bereg, Holger H Hoos, and Anne Condon. RNA STRAND: The RNA secondary structure and statistical analysis database. *BMC bioinformatics*, 9(1):340, 2008.
- Mirela Andronescu, Anne Condon, Holger H Hoos, David H Mathews, and Kevin P Murphy. Computational approaches for RNA energy parameter estimation. *RNA*, 16(12):2304–2318, 2010.
- Liron Argaman and Shoshy Altuvia. fhla repression by oxys RNA: kissing complex formation at two sites results in a stable antisense-target RNA complex. *Journal of molecular biology*, 300(5):1101–1112, 2000.
- David P Bartel. MicroRNAs: genomics, biogenesis, mechanism, and function. *cell*, 116(2):281–297, 2004.
- Chase L Beisel and Gisela Storz. The base-pairing RNA spot 42 participates in a multioutput feedforward loop to help enact catabolite repression in escherichia coli. *Molecular cell*, 41(3):286–297, 2011.
- Isaac Bentwich. Prediction and validation of micrnas and their targets. *FEBS letters*, 579(26):5904–5910, 2005.
- Philip N Borer, Barbara Dengler, Ignacio Tinoco Jr, and Olke C Uhlenbeck. Stability of ribonucleic acid double-stranded helices. *Journal of molecular biology*, 86(4):843–853, 1974.
- Julius Brennecke, Alexander Stark, Robert B Russell, and Stephen M Cohen. Principles of microrna–target recognition. *PLoS biology*, 3(3), 2005.
- Anke Busch, Andreas S Richter, and Rolf Backofen. IntaRNA: efficient prediction of bacterial srna targets incorporating target site accessibility and seed regions. *Bioinformatics*, 24(24):2849–2856, 2008.
- Hamidreza Chitsaz, Rolf Backofen, and S Cenk Sahinalp. biRNA: Fast RNA-RNA binding sites prediction. *International Workshop on Algorithms in Bioinformatics*, pages 25–36, 2009a.

- Hamidreza Chitsaz, Raheleh Salari, S Cenk Sahinalp, and Rolf Backofen. A partition function algorithm for interacting nucleic acid strands. *Bioinformatics*, 25(12):i365–i373, 2009b.
- Howard DeVoe and Ignacio Tinoco. The stability of helical polynucleotides: base contributions. *Journal of molecular biology*, 4:500–17, 1962.
- Robert M Dirks and Niles A Pierce. A partition function algorithm for nucleic acid secondary structure including pseudoknots. *Journal of computational chemistry*, 24(13):1664–1677, 2003.
- Robert M Dirks, Justin S Bois, Joseph M Schaeffer, Erik Winfree, and Niles A Pierce. Thermodynamic analysis of interacting nucleic acid strands. *SIAM review*, 49(1):65–88, 2007.
- Tsukasa Fukunaga and Michiaki Hamada. Riblast: an ultrafast rna–rna interaction prediction system based on a seed-and-extension approach. *Bioinformatics*, 33(17):2666–2674, 2017.
- Rick Gelhausen. Constrained RNA-RNA interaction prediction, 2018.
- Michael Ibba and Dieter Söll. Aminoacyl-tRNA synthesis. *Annual review of biochemistry*, 69(1):617–650, 2000.
- Edwin T Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620, 1957.
- Tamás Kiss. Small nucleolar RNAs: an abundant group of noncoding RNAs with diverse cellular functions. *Cell*, 109(2):145–148, 2002.
- Martin Mann, Patrick R Wright, and Rolf Backofen. IntaRNA 2.0: enhanced and customizable prediction of RNA–RNA interactions. *Nucleic acids research*, 45(W1):W435–W439, 2017.
- David H Mathews, Jeffrey Sabina, Michael Zuker, and Douglas H Turner. Expanded sequence dependence of thermodynamic parameters improves prediction of RNA secondary structure. *Journal of molecular biology*, 288(5):911–940, 1999.
- David H Mathews, Matthew D Disney, Jessica L Childs, Susan J Schroeder, Michael Zuker, and Douglas H Turner. Incorporating chemical modification constraints into a dynamic programming algorithm for prediction of rna secondary structure. *Proceedings of the National Academy of Sciences*, 101(19):7287–7292, 2004.
- John S McCaskill. The equilibrium partition function and base pair binding probabilities for RNA secondary structure. *Biopolymers: Original Research on Biomolecules*, 29(6-7):1105–1119, 1990.
- Barmak Modrek and Christopher Lee. A genomic view of alternative splicing. *Nature genetics*, 30(1):13, 2002.
- Ulrike Mückstein, Hakim Tafer, Jörg Hackermüller, Stephan H Bernhart, Peter F Stadler, and Ivo L Hofacker. Thermodynamics of RNA–RNA binding. *Bioinformatics*, 22(10):1177–1182, 2006.
- Saul B Needleman and Christian D Wunsch. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of molecular biology*, 48(3):443–453, 1970.
- Dmitri D Pervouchine. Iris: intermolecular RNA interaction search. *Genome Informatics*, 15(2):92–101, 2004.

- Sarah C Pulvermacher, Lorraine T Stauffer, and George V Stauffer. The role of the small regulatory RNA *gcvB* in GcvB/mRNA posttranscriptional regulation of *oppA* and *dppA* in *escherichia coli*. *FEMS microbiology letters*, 281(1):42–50, 2008.
- Martin Raden, Mostafa Mahmoud Mohamed, Syed Mohsin Ali, and Rolf Backofen. Interactive implementations of thermodynamics-based RNA structure and RNA–RNA interaction prediction approaches for example-driven teaching. *PLoS computational biology*, 14(8):e1006341, 2018.
- Maria Selmer, Christine M Dunham, Frank V Murphy, Albert Weixlbaumer, Sabine Petry, Ann C Kelley, John R Weir, and Venki Ramakrishnan. Structure of the 70s ribosome complexed with mRNA and tRNA. *Science*, 313(5795):1935–1942, 2006.
- Brian Tjaden, Sarah S Goodwin, Jason A Opdyke, Maude Guillier, Daniel X Fu, Susan Gottesman, and Gisela Storz. Target prediction for small, noncoding RNAs in bacteria. *Nucleic acids research*, 34(9):2791–2802, 2006.
- Sinan Uğur Umu and Paul P Gardner. A comprehensive benchmark of RNA–RNA interaction prediction tools for all domains of life. *Bioinformatics*, 33(7):988–996, 2017.
- Mark L Urbanowski, Lorraine T Stauffer, and George V Stauffer. The *gcvB* gene encodes a small untranslated RNA involved in expression of the dipeptide and oligopeptide transport systems in *escherichia coli*. *Molecular microbiology*, 37(4):856–868, 2000.
- Patrick R Wright, Andreas S Richter, Kai Papenfort, Martin Mann, Jörg Vogel, Wolfgang R Hess, Rolf Backofen, and Jens Georg. Comparative genomics boosts target prediction for bacterial small RNAs. *Proceedings of the National Academy of Sciences*, 110(37):E3487–E3496, 2013.
- Patrick R Wright, Martin Mann, and Rolf Backofen. Structure and interaction prediction in prokaryotic RNA biology. *Microbiol Spectrum*, 6(2):10–1128, 2018.
- Xiaojun Xu and Shi-Jie Chen. Physics-based RNA structure prediction. *Biophysics reports*, 1(1):2–13, 2015.
- Zhigang Zhang, Gia-Phong Vu, Hao Gong, Chuan Xia, Yuan-Chuan Chen, Fenyong Liu, Jianguo Wu, and Sangwei Lu. Engineered external guide sequences are highly effective in inhibiting gene expression and replication of hepatitis b virus in cultured cells. *PloS one*, 8(6), 2013.
- Michael Zuker and Patrick Stiegler. Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information. *Nucleic acids research*, 9(1):133–148, 1981.

Appendices

Appendix A

RNA Sequences

| RNA | Sequence |
|---------|--|
| OxyS | GAAACGGAGCGGCACCTCTTTTAACCCTTGAAGTCACTGCCCCGTTTCGAGAGTTTCTCAACTCGAATAACTAAAGC CAACGTGAACCTTTTGCGGATCTCCAGGATCCGC |
| fh1A | AGTTAGTCAATGACCTTTTGCACCGCTTTGCGGTGCTTTCCTGGAAGAACAAAATGTCATATACACCGATGAGTGGA TCTCGGACAAACAAGGGTTGTTGACATCACTCGGACA |
| Spot42 | GUAGGGUACAGAGGUAAGAUGUUCUAUCUUCAGACCCUUUACUUCACGUAAUCGGAUUUGGCUGAAUAAUUUAGC CGCCCCAGUCAGUAAUGACUGGGGCGUUUUUA |
| sthA | GGGATCAATTGGCTTACCCGCGATAAAATGTTACCATTCTGTTGCTTTTATGTATAAGAACAGGTAAGCCCTACCA TGCCACATTCTTACGATTACGATGCCATAGTAATAGGTTCCGGCCCCGGCGGCCGAAGGCGCTGCAATGGGCCTG |
| gcvB | TTCTTGAGCCGGAACGAAAAGTTTATCGGAATGCGTGTCTGTATGGGCTTTTGGCTTACGGTTGTGATGTTGTGT TGTTGTGTTTGAATTGGTCTGCGATTTCAGACCACGGTAGCGAGACTACCCTTTTCACTTCCTGTACATTTACCC TGTCTGTC |
| oppA | GACAGCAGAAAGUCUCCGAGCCUGUGCAGGGUCCCAAUCCGGGAUUACACAUGCUGGUAAUACCAGUAAUUUAAA UGAGGGAGUCCAAAAAACAAUGACCAACAUCACCAAGAGAAGUUUAGUAGCAGCUGGCGUUCUGGCUGCGCUAAUG GCAGGGA |
| DicF | TTTCTGCTGACGTTTGGCGGTATCAGTTTTACTCCGTGACTGCTCTGCCGCCCTTTTAAAGTGAATTTT |
| ftsZ | AAAAGAGTTTTAATTTTATGAGGCCGACGATGATTACGGGCTCAGGCGACAGGCACAAATCGGAGAGAACTATG TTTGAAACCAATGGAACCTTACCAATGACGCGGTGATTAAAGTCATCGGCGTGGCGGCGGCGGCGGTAATGCTGTTG AACACATGGTGGCGGAGCGCATTGAAGGTGTTGAATTCCTCGCGGTAAATACCGATGCACAAGCGCTGCGTAAAA |
| EGS | AACCTTGCTCGTTATCGCTGGATGTGTCTGCGGCGTTTTATCATCTTCCTCTTCATCCTGTGCTATGCCTCATC TTCTTGTTGGTTCTTCTGGACTATCAAGGTAAGTTGCCCCGTTTGCTCTTAAT |
| S-SER | GTTAACGATGAGGGTGCGGTCTCCGCGCGCAGGTTCAAATCCTGCTAGCAGCATTT |
| S-SER-C | GTTAACGATGAGGGTGCGGTCTCCGCGCGCAGGAAGAAATCCTGCTAGCAGCATTT |
| S-C386 | GTTAACGATGAGGGACCTCACCGGTCTGTTCCGATTTCGACTAGCAGCATTT |
| S-C386C | GTTAACGATGAGGGACCTCACCGGTCTGGAAGGATTTCGACTAGCAGCATTT |

Table A.1: Table of RNAs with their corresponding sequence as used in this thesis.