

# **RNA STRUCTURE ANALYSIS VIA THE RIGID BLOCK MODEL**

by  
**MAURICIO ESGUERRA NEIRA**

A dissertation submitted to the  
Graduate School—New Brunswick  
Rutgers, The State University of New Jersey  
in partial fulfillment of the requirements  
for the degree of  
Doctor of Philosophy  
Graduate Program in Chemistry and Chemical Biology

Written under the direction of  
Wilma K. Olson  
and approved by

---

---

---

---

New Brunswick, New Jersey  
May, 2010

## **ABSTRACT OF THE DISSERTATION**

### **RNA Structure Analysis via the Rigid Block Model**

**by Mauricio Esguerra Neira**

**Dissertation Director: Wilma K. Olson**

RNA structure is at the forefront of our understanding of the origin of life, and the mechanisms of life regulation and control. RNA plays a primordial role in some viruses. Our knowledge of the importance of RNA in cellular regulation is relatively new, and this knowledge, along with the detailed structural elucidation of the transcription machine, the ribosome, has propelled interest in understanding RNA to a level which starts to closely resemble that given to proteins and DNA.

In the process of progressively understanding the landscape of functionality of such a complex polymer as RNA, one practical task left to the structural chemist is to understand the details of how structure relates to large-scale polymer processes. With this in mind the fundamental problems which fuel the work described in this thesis are those of the conformations which RNA's assume in nature, and the aim to understand how RNA folds.

The RNA folding problem can be understood as a mechanical problem. Therefore effort to determine its solution are not foreign to the use of statistical mechanical methods combined with detailed knowledge of atomic level structure. Such methodology is mainly used in this work in a long term effort to understand the intrinsic structural features of RNA, and how they might relate to its folding.

*As a thing among things, each thing is equally insignificant; as a world each one equally significant.*

*If I have been contemplating the stove, and then am told; but now all you know is the stove, my result does indeed sound trivial. For this represents the matter as if I had studied the stove as one among the many, many things in the world. But if I was contemplating the stove, it was my world, and everything else colorless by contrast with it ...*

*For it is equally possible to take the bare present image as the worthless momentary picture in the whole temporal world, and as the true world among shadows.*

**Ludwig Wittgenstein**

*As a molecule among molecules, each molecule is equally insignificant; as a world each one equally significant.*

*If I have been contemplating RNA, and then am told; but now all you know is RNA, my result does indeed sound trivial. For this represents the matter as if I had studied RNA as one among the many, many molecules in the world. But if I was contemplating RNA, it was my world, and everything else colorless by contrast with it ...*

*For it is equally possible to take the bare present image as the worthless momentary picture in the whole temporal world, and as the true world among shadows.*

**Anonymous Chemist**

## **Acknowledgements**

I would first like to give a special thanks to Dr. Yurong Xin, whose patience, help, and collaboration since the very beginning of my joining of the Olson lab have been fundamental for the development of this work. I would like to thank Dr. Olson's extreme patience, and room for freedom on carrying out this research. Finally I thank all colleagues at the Olson lab.

# Table of Contents

Abstract . . . . .	ii
Acknowledgements . . . . .	iv
<b>List of Tables . . . . .</b>	<b>vii</b>
<b>List of Figures . . . . .</b>	<b>viii</b>
<b>1. Introduction . . . . .</b>	<b>1</b>
1.1. RNA . . . . .	1
1.1.1. RNA folding . . . . .	1
1.2. Is RNA folding a hard or easy problem? . . . . .	2
1.3. Experimental folding techniques . . . . .	2
1.4. RNA simulations . . . . .	3
1.4.1. Local nucleotide interactions . . . . .	3
1.4.2. RNA secondary structure algorithms and the lack of tertiary ones . . . . .	4
1.4.3. RNA overall fold . . . . .	4
1.4.4. RNA motifs . . . . .	5
1.5. Overview . . . . .	6
References . . . . .	7
<b>2. RNA Base Steps . . . . .</b>	<b>13</b>
2.1. Consensus Clustering of Single Stranded Base Step Parameters . . . . .	14
2.1.1. Combining Fourier Averaging Results and Clustering Analysis . . . . .	15
2.1.2. Partitional Clustering for Rigid Body Parameters . . . . .	16
2.1.3. Hierarchical Clustering for Rigid Body Parameters . . . . .	21
2.2. RNA Conformations . . . . .	21
References . . . . .	27
<b>3. RNA Base-Pairing . . . . .</b>	<b>29</b>
3.1. Canonical and Noncanonical Base-pairs, Methods Paper . . . . .	29
3.2. Clustering of Yurong's Classification . . . . .	29
<b>4. RNA Base Pair Steps . . . . .</b>	<b>30</b>
4.1. Analysis (Albany Poster) and Django Webserver . . . . .	30
4.2. Persistence Length vs. Hagerman . . . . .	30
4.3. AMBER: Persistence Length of Base-Pair Step Patterns . . . . .	30
<b>5. RNA Motifs . . . . .</b>	<b>31</b>
5.1. GNRA tetraloop . . . . .	31
5.1.1. 3DNA-Parser . . . . .	31
5.1.2. Overlap Scores . . . . .	31
5.2. Triplets on RNA (comparison to Laing et al.) . . . . .	32

References . . . . .	35
<b>6. RNA Helical Regions and Graph Theory . . . . .</b>	<b>36</b>
<b>Appendix A. Clustering Analysis (CA) . . . . .</b>	<b>37</b>
A.1. Hierarchical methods . . . . .	37
<b>Appendix B. Dimension Reduction . . . . .</b>	<b>40</b>
B.1. Principal Component Analysis . . . . .	40
References . . . . .	41
<b>Appendix C. Figure Supplements . . . . .</b>	<b>42</b>
C.1. Chapter2 . . . . .	42
Curriculum Vitae . . . . .	44

## List of Tables

2.1. Some large RNA structures (>300 bases) elucidated in the last decade. . . . .	15
2.2. Residue numbers for base-steps with RMSD values less than 15 between the reference base-step vectors from the four groups of non-A-type RNA dinucleotide conformations and all base-step vectors found in the 23S strand of <i>Haloarcula marismortui</i> large ribo- somal subunit. . . . .	19
2.3. Base step torsion angles for the different known RNA conformations. . . . .	21
2.4. Base step parameters for the different known RNA conformations. Notice that the base step parameters are for single bases rather than base-pairs. . . . .	21
A.1. Example of structures, considered as bidimensional vectors, to be clustered using the average linkage method and the Manhattan distance. . . . .	38

## List of Figures

1.1.	Separation of secondary and tertiary interaction in RNA [27]. Double helical secondary structure represented by individual cylinders and tertiary interactions by association of cylinders. . . . .	2
1.2.	Ribbon-coil schematic illustrating the fold and intermolecular units of a dimer of prealbumin, or transthyretin, taken from Richardson <i>et al.</i> [73] . . . . .	5
1.3.	<i>Haloharcula marismortui</i> 's large ribosomal subunit (left) and hammerhead ribozyme (right). The figures were taken directly from the NDB web pages, and show a ribbon representation of the phosphate backbone, and a block representation for the nucleotide bases. From the figures it's clear that, whereas the ribozyme fold can be clearly understood with this representation, the ribosome fold cannot. . . . .	5
2.1.	<b>Right:</b> Total number of RNA bases added to the PDB database between 2000 and 2010 (Exponential fit line in blue). <b>Left:</b> Total number of RNA structures solved yearly by X-Ray crystallography between 2000 and 2010 (Exponential fit line in red). . . . .	13
2.2.	Frequency of nucleotide bases in RNA molecules found in the PDB classified by the size of RNA molecules. We define the size as the total number of nucleotide bases present per molecule. . . . .	14
2.3.	Figure taken from Richardson et al. [11] where the blue and green dots in a) mean very accurate van der Waals distances, and in b) the red and orange dots mean steric clashes, that is, distances outside the acceptable van der Waals range. . . . .	15
2.4.	Dendrogram showing the results of consensus clustering of 20 non-Atype rRNA dinucleotides according to their hexadimensional base-step parameter vectors. . . . .	17
2.5.	rRNA dinucleotide structures organized by clusters obtained from consensus clustering of their hexadimensional base-step parameter vectors. . . . .	18
2.6.	Sum of all within clusters sum of squares against number of clusters. . . . .	20
2.7.	Average silhouette width against number of clusters. . . . .	20
2.8.	Cluster dissimilarities for the twelve hierarchical trees obtained from clustering of the six-dimensional base-step parameters obtained from the large subunit of the ribosome (PDB-ID:1jj2) . . . . .	22
2.9.	K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the <i>Hartigan-Wong</i> algorithm. The number of partitions is <b>2</b> . The upper diagonal matrix displays the values of the linear correlation coefficient $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal. . . . .	23
2.10.	K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the <i>Lloyd</i> algorithm. The number of partitions is <b>2</b> . The upper diagonal matrix displays the values of the linear correlation coefficient $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal. . . . .	24
2.11.	K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the <i>Forgy</i> algorithm. The number of partitions is <b>2</b> . The upper diagonal matrix displays the values of the linear correlation coefficient $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal. . . . .	25



2.12.	K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the <i>McQueen</i> algorithm. The number of partitions is <b>2</b> . The upper diagonal matrix displays the values of the linear correlation coefficient $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal. . . . .	26
5.1.	GNRA Tetraloop from <i>Thermus Thermophilus</i> 23S Ribosomal RNA PDB-ID:1ffk. . . . .	32
5.2.	Normalized histograms showing the distribution of overlap values in the 23S subunit or <i>Thermus Thermophilus</i> rRNA, PDB-ID:1jjk. In histogram (a) all values are included, but in histogram (b) only values greater than zero are included. Notice the high preponderance of zero values, exactly 897 out of a total of 2705. . . . .	33
5.3.	Dendrogram for consensus clustering of overlap scores in the ribosome. Zero values filtered out and remaining data normalized. . . . .	34
A.1.	Clustering tree for 5 bidimensional vectors using the Manhattan distance definition and the average linkage clustering method. . . . .	39
C.1.	Non A-RNA Type base steps centered on the standard reference frame of Adenine. Top view with the Minor Groove side of Adenine pointing down the page and the Major Groove pointing up. . . . .	43

# Chapter 1

## Introduction

### 1.1 RNA

RNA plays a primordial role in life, and perhaps also in the early history of its origins [1, 2, 3, 4]. In Biology RNA is a central player in the transcription and translation steps of what is known as its central dogma, i.e., DNA makes RNA (via transcription) and RNA makes protein (during translation). In the last decade of the twentieth century Fire and Mello [5] found that RNA also plays a role previously thought to be the job of proteins. That is, RNA can regulate translation using non-coding RNA's (ncRNA's). Another fundamental discovery about RNA came in 2000 with the elucidation at atomic level detail of a non-coding RNA, the ribosome [6, 7, 8].

Since its very beginnings, structural understanding of RNA has proven to be a very complex problem. It was not until 1956, three years after the famous Nature triad of papers by Watson and Crick, Wilkins, Stoke, and Wilson, and Franklin and Gosling [9] on the double stranded structure of DNA, that Alex Rich and David Davies were able to produce double stranded RNA from polyriboadenylic acid (poly-rA) and polyribouridylic acid (poly-rU) to produce a neatly diffracting X-ray pattern typical of a double-helical structure. It was not until 1965 that Robert Holley was able to obtain the complete sequence of yeast Alanine tRNA, and also its secondary structure from cleavage of the whole structure into smaller fragments, and it was only in 1973, that the first complex, but small, tRNA structure, was solved at full atomic detail. Fifty seven years have passed since the description of the double-helical structure of DNA, but still RNA faces more challenges with the possibility of finding a whole new zoo of non-coding RNA structures [10], and the possibility of new engineered ones [11].

#### 1.1.1 RNA folding

The first high-resolution X-ray structure of RNA larger than a dinucleotide was that of yeast tRNA<sup>Phe</sup> at 3 Å in 1974 [12, 13, 14]. Thirty six years later there are two orders of magnitude more structural information about RNA [15], and new information from non-coding RNA's is expected [10]. This fact and the discovery of ribozymes [16, 17], which are catalytic RNA molecules, has renewed interest in solving the RNA folding problem, that is, starting from the primary sequence, finding in an automated<sup>i</sup> way the native three-dimensional structure of an RNA molecule and the folding pathway that it follows. The RNA folding problem is usually seen as analogous to the protein folding problem, due both to the discovery of the enzymatic behavior of RNA [16, 17] and the complicated folding of large RNA molecules [21]. To take advantage of this analogy, a unified conceptual framework for describing RNA and protein folding, called the kinetic partitioning mechanism (KPM), has been developed by Thirumalai and Hyeon [22]. This and other methods are based on defining an adequate partition function for describing the correct conformational ensemble of folded, partially folded, and unfolded structures [23, 24, 25] of either protein or RNA.

---

<sup>i</sup>The term automated is used here to mean a theoretical model of tertiary folding, which could use experimental measures of secondary structure association in the same way that the traditional secondary structure folding model [18, 19] uses the Tinoco-Uhlenbeck dinucleotide postulate [20] to find total free energies.

## 1.2 Is RNA folding a hard or easy problem?

There are two trains of thought regarding the mechanism of RNA folding. One states that RNA folding is less complex than protein folding [26] because RNA is made up of a four letter alphabet of similar nucleotide units instead of a 20 letter alphabet of dissimilar amino acids. Therefore the number of possible sequential combinations is smaller. It is also well known that secondary and tertiary interactions can be separated in the case of RNA by the absence or presence of  $Mg^{2+}$  [27] (see Figure 1.1), and that the secondary structure motifs of RNA are more limited in number than those of protein, whereas secondary and tertiary elements are not as easily separable in proteins. The other point of view says

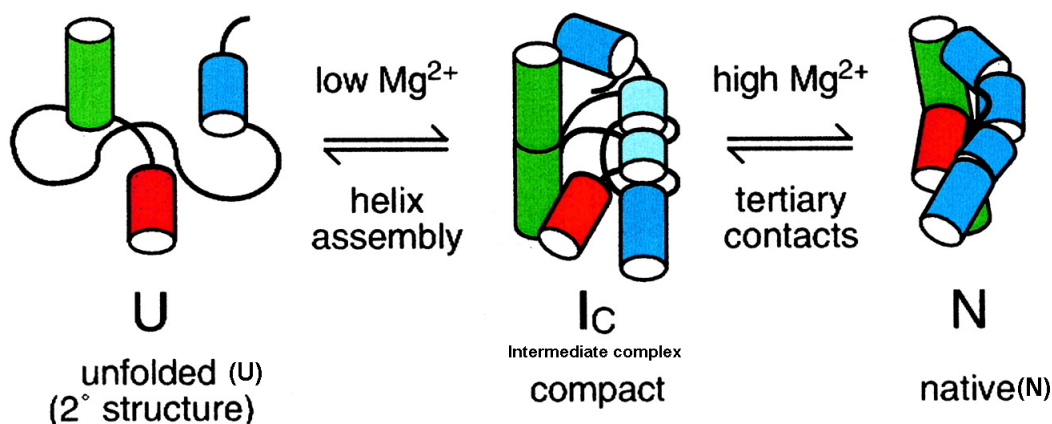


Figure 1.1: Separation of secondary and tertiary interaction in RNA [27]. Double helical secondary structure represented by individual cylinders and tertiary interactions by association of cylinders.

that RNA folding can be at least as complex as protein folding [28, 29] since there is no such thing as hydrophobic burial of regions of RNA as in the case of proteins. Instead, the electrostatic problem of having a complex charged backbone must be dealt with in the case of RNA. For instance, the interactions of the RNA polyanionic backbone with water and cations [30] are not easily simulated with explicit solvent models as can be done for proteins. The aforementioned interactions of RNA need to be modeled implicitly, and must aim to describe long dynamic processes of the order of seconds to minutes, in contrast to the typical time scales of tens of microseconds associated with protein folding. Although secondary and tertiary structure can be separated experimentally, there have been few theoretical efforts to account for the folding of RNA from a random sequence of nucleotides into secondary structures and tertiary structures. What little is known has been investigated at low resolution. Professor Stephen Harvey and associates have simulated yeast tRNA<sup>Phe</sup>, [31] and the assembly of the 30S subunit of the ribosome [32] at various levels of detail, initially using only one pseudoatom per helical region, and later one pseudoatom per nucleotide. Recently Major's group at Montreal has proposed a pipeline of two computer algorithms [33], one makes secondary structure predictions, and the other assembles 3D structures based on the best scoring secondary structures. By contrast, in the case of proteins many groups have simulated the transition from secondary to tertiary structure, including some calculations which account for the strong coupling of secondary and tertiary structure [34, 35, 36]. This type of work is often referred to as protein structural topology and there is no counterpart for RNA.

## 1.3 Experimental folding techniques

Traditionally RNA folding and unfolding have been followed calorimetrically and spectroscopically as a function of temperature and cation concentration [37, 38]. While this approach works well for studying two-state folders, *i.e.*, structures which populate only two states (native and melted), in general RNA's

are not two-state folders. RNA seems to go through a rugged free energy landscape of conformations in the process of folding [39]. The experimental solution to this problem is offered by single molecule techniques like fluorescence resonance energy transfer (FRET) and mechanical micromanipulation, in which the ends of RNA are attached to micron sized beads which are then pulled apart and monitored with a laser light trap [40, 41, 42, 43]. In the case of single molecule force-induced unfolding, state transitions often occur under non-equilibrium conditions, thereby making it difficult to extract equilibrium information from the data. Recently Bustamante, Tinoco, and associates have shown that using the Crooks fluctuation theorem [44], one can deal with such cases and extract RNA folding free energies from single molecule experiments [45].

## 1.4 RNA simulations

Network and molecular mechanics-molecular dynamics (MM-MD) methods provide useful information relevant to the RNA folding-unfolding problem, especially for describing fluctuations away from the native conformation. Gaussian network models [46, 47, 48] which treat RNA at less than atomic detail have been used to describe the motions of large RNA structures like the ribosome. Examples of the predicted normal modes of motion of the ribosome can be seen at: <http://ribosome.bb.iastate.edu/70SnKmode>. Using MM, Sanbonmatsu and coworkers obtained a static atomic model of the 70S ribosome structure through homology modeling [49]. Tung and associates used this structure for an all-atom MD simulation of the movement of tRNA into a fluctuating ribosome [50]. This type of simulation might be useful in a reverse-folding approach to the RNA folding problem. To the best of our knowledge, such calculations haven't as yet been done for RNA.

### 1.4.1 Local nucleotide interactions

The molecular interactions which rule RNA structures at the nucleic acid base level, *i.e.*, local level, are hydrogen bonding and stacking interactions. The former are related to base pairing and the latter, in most cases, to nucleotide steps. These interactions can be explored theoretically at various levels. At the highest level are ab-initio quantum mechanical calculations which are still too expensive for systems as large as hundreds of atoms. Such calculations, nevertheless, can tell a great deal about local electronic behavior. For example, Hobza and collaborators have found that the stacking interaction of free nucleotide bases is determined by dispersion attraction, short-range exchange repulsion, and electrostatic interaction. No specific  $\pi - \pi$  interactions are found from electron correlated ab-initio calculations [51, 52]. This is why force field methods have been so successful in the study of nucleic acids, since the empirical potentials used in such studies mimic well the quantum mechanically obtained energy profiles [49, 53]. A currently debated ab-initio finding is whether small fluctuations in the configurations of neighboring base pairs (dimers) are iso-energetic or not. Recent calculations of Sponer and Hobza [54] seem to contradict their older publications [53, 55], in which the stacking energies were reported to be relatively insensitive to dimer conformation. The new results use the so-called "coupled cluster singles doubles with triple electron excitations" CCSD(T) method, to account for electron correlation. Using this electron correlation energy correction, the stacking energy differences between dimer conformations turn out to be considerably higher than previously reported.

Single and double strand stacking free energies can be obtained calorimetrically. The most popular method used for obtaining such quantities is differential scanning calorimetry (DSC) [56]. These measurements show favorable dinucleotide stacking free energies as large as -3.6 kcal/mol for double strand stacking. Experimentally, the magnitudes of these interactions are found to be sequence dependent [37]. In fact, the stacking free energies for some sequences<sup>ii</sup> are found to be negligible. Thus there

---

<sup>ii</sup>Unpaired terminal nucleotides UC/A UU/A at 1M NaCl.

may be no accountable stacking interaction at all for some sequences.

Besides taking into account the effects of stacking and hydrogen bonding, it is important to think at the same time about the polyelectrolyte nature of the RNA backbone. Manning's counterion condensation theory [57, 58] provides a simple and quantitative picture of the interactions of the double helical nucleic acid polyanion with its counterions, although it does not take into account the discrete nature of charge [37] or the folding of RNA. Poisson-Boltzmann theory offers a more detailed picture of the behavior of charged macroions in solution [59].

The local conformational space of RNA has been studied using a large set of available RNA structures from the Nucleic Acid Database (NDB) [60]. The torsion angles of the nucleotide steps have been clustered in the parameter space using different techniques [61, 62]. The root-mean-square deviations (RMSD) of the distances between closely spaced atoms in the phosphates, sugars, and bases, have also been clustered [63]. The latter studies are aimed at finding the common nucleotide base steps and base-pair building blocks which are given the name of RNA doublets. Recently, the RNA Ontology Consortium (ROC) has proposed a consensus set of RNA dinucleotide conformers integrating the work of various groups [64].

### 1.4.2 RNA secondary structure algorithms and the lack of tertiary ones

From secondary structure prediction algorithms like Zuker's *mfold* program [65], Hofacker's Vienna RNA package [19], or Mathews Dynaling [66], one obtains a large ensemble of secondary structure graphs. These graphs can be analyzed with graph theory to produce a partition function describing a full arrangement of contacts for the total number of possible secondary structures making possible a "relation of microscopic conformations to macroscopic properties" [67]. So far this type of model has not been generalized to take into account tertiary structural features, *i.e.*, interhelical interactions of RNA. In the last two to three years a boom in prediction of small ( $\approx 200$  nucleotides) RNA 3D structures has started. Basically three types of approaches are being followed. One is that of using a coarse grained model assigning a potential function to it, followed by a minimization procedure, and then a molecular mechanics (MM) all atom refinement [68, 69, 70]. Another starts from predicted secondary structures and assumes their helical regions adopt the A-form conformation, then mechanically thrusts residues as rigid bodies in the remaining non-helical regions, and finally carry out an MM optimization [71]. Finally, a pipeline between secondary structure prediction, and tertiary structure assembly is proposed. This pipeline uses as bridging concept between 2D and 3D structure, the graph theoretical definition of a minimum cycle basis, which for the case of nucleic acids is renamed by Major's group as Nucleic Cyclic Motifs (NCM) [33].

### 1.4.3 RNA overall fold

Whereas in the case of proteins one can describe the overall fold from the arrangement of secondary structure motifs, *i.e.*, using the helix-ribbon-coil images developed by Jane Richardson [72] (see Figure 1.2), there is still no comparable description of the overall fold of RNA. A ribbon representation of the sugar phosphate backbone helps to understand the folding of small RNA's, but in the case of the ribosome this type of representation is not sufficient, see Figure 1.3.

One can envision that a thorough investigation of the parameter space of translational and rotational degrees of freedom of the helical regions of RNA could give clues as to how we might see an overall fold in RNA structures.

In the case of proteins the SCOP (Structural Classification of Proteins) database [74], classifies proteins, among other classifications, according to recurrent arrangements of secondary structure, that is, folds. The SCOR (Structural Classification of RNA) database [75, 76], aims to provide a similar

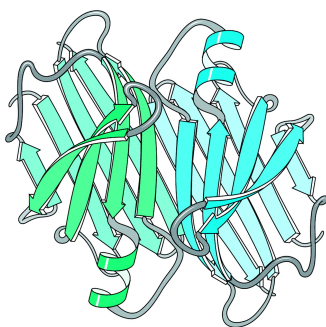


Figure 1.2: Ribbon-coil schematic illustrating the fold and intermolecular units of a dimer of prealbumin, or transthyretin, taken from Richardson *et al.* [73]

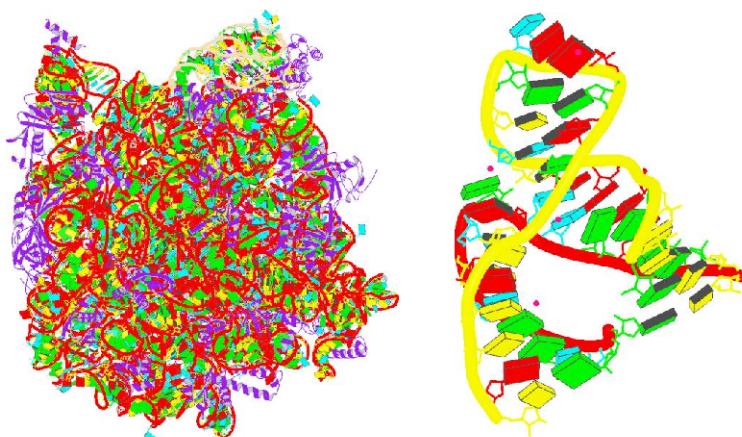


Figure 1.3: *Haloharcula marismortui*'s large ribosomal subunit (left) and hammerhead ribozyme (right). The figures were taken directly from the NDB web pages, and show a ribbon representation of the phosphate backbone, and a block representation for the nucleotide bases. From the figures it's clear that, whereas the ribozyme fold can be clearly understood with this representation, the ribosome fold cannot.

classification to that obtained for proteins, but using RNA motifs instead. This classification focuses on the local folding of small pieces of RNA and cannot describe the overall fold.

#### 1.4.4 RNA motifs

First, a word of caution must be given to the reader. The term “*RNA motif*” alone is used in the literature to describe three different levels of RNA organization, that is, RNA **sequence** motifs, RNA **secondary structure** motifs, or RNA **3D structure** motifs. We start by making such distinction as it is not always clearly mentioned in RNA literature, generating a great deal of confusion and bibliographical search frustration for the beginner. The kind of RNA motifs we are dealing with in this thesis are those of the third kind, that is, RNA **3D structure** motifs which we'll address from now on simply as RNA motifs. Yet another source of confusion in understanding RNA motifs is the lack of a unique definition. Three popular and somewhat recent definitions are:

- RNA motifs are “*Conserved structural subunits that make up the secondary structures of RNAs.*”[77]
- RNA motifs are “*Ordered stacked arrays of non-Watson-Crick base pairs that form distinct folds on the phosphodiester backbones of RNA strands.*”[78]

- “An RNA Motif is a discrete sequence or combination of base juxtapositions found in naturally occurring RNA's in unexpectedly high abundance.”[79]

From our point of view RNA motifs are to be understood as peculiar sets of geometrical (in the rigid block sense) arrangements in three dimensional space.

Even though there is no unique definition, we can think of three practical tasks regarding RNA motifs. That is, given an RNA 3D structure automatically identify [80, 81, 82], describe [83, 84, 85, 86, 87], and find new [88, 89, 92, 90, 81] motifs.

## 1.5 Overview

Keeping always in mind the greater scope of the RNA folding problem, this thesis addresses various issues of RNA structural understanding using RNA crystallographic data from the Protein Data Bank (PDB). Such data has been analyzed statistically along with the use of a very rigorous rigid body formalism. In Chapter 2 the consensus clustering technique is used to classify RNA base-step parameters of non-ARNA conformations, and the resulting groups are localized and understood in the context of rRNA. Chapter 3 reconsiders previous work carried out by Dr. Yurong Xin at the Olson's lab, on classification of RNA base-pairs by resorting again to clustering analysis techniques, and database mining of the WWW available Base Pair Structures (BPS) database. In Chapter 4 we explore, using statistical analysis, the data available on RNA helical regions, and use it to compute the persistence length of double stranded RNA's and compare it to experimental results. In Chapter 5 we provide a new python software, pyRNAmotifs which interfaces with 3DNA to do a rigorous search of existing and perhaps new RNA motifs, and finally in Chapter 6 we propose the measurement and classification of RNA structures using a new graph theoretical index named folding index, based on a helical region "view" of RNA's, which is clearly concordant with the emerging necessity of new metrics beyond RMSD for structural understanding.

## References

- [1] Woese, C. (1967) The Genetic Code, the Molecular Basis for Genetic Expression, Harper and Row, .
- [2] Crick, F. (1968) The Origin of the Genetic Code. *Journal of Molecular Biology*, **38**, 367–379.
- [3] Orgel, L. (1968) Evolution of the Genetic Apparatus. *Journal of Molecular Biology*, **38**, 381–393.
- [4] Orgel, L. E. (2004) Prebiotic Chemistry and the Origin of the RNA World. *Critical Reviews in Biochemistry and Molecular Biology*, **39**, 99–123.
- [5] Fire, A., Xu, S., Montgomery, M. K., Kostas, S. A., Driver, S. E., and Mello, C. C. (Feb, 1998) Potent and Specific Genetic Interference by Double-Stranded RNA in *Caenorhabditis Elegans*. *Nature*, **391**(6669), 806–811.
- [6] Schluenzen, F., Tocilj, A., Zarivach, R., Harms, J., Gluehmann, M., Janell, D., Bashan, A., Bartels, H., Agmon, I., Franceschi, F., and Yonath, A. (2000) Structure of Functionally Activated Small Ribosomal Subunit at 3.3 Å Resolution. *Cell*, **102**, 615–623.
- [7] Ban, N., Nissen, P., Hansen, J., Moore, P. B., and Steitz, T. A. (August, 2000) The Complete Atomic Structure of the Large Ribosomal Subunit at 2.4 Å Resolution. *Science*, **289**(5481), 905–920.
- [8] Wimberly, B. T., Brodersen, D. E., Clemons, W. M., Morgan-Warren, R. J., Carter, A. P., Vonnrhein, C., Hartschk, T., and Ramakrishnan, V. (2000) Structure of the 30S Ribosomal Subunit. *Nature*, **407**, 327–339.
- [9] Watson, J. D. and Crick, F. H. (Apr, 1953) Molecular Structure of Nucleic Acids; A Structure for Deoxyribose Nucleic Acid. *Nature*, **171**(4356), 737–738.
- [10] Weinberg, Z., Perreault, J., Meyer, M. M., and Breaker, R. R. (2009) Exceptional Structured Non-coding RNAs Revealed by Bacterial Metagenome Analysis. *Nature*, **462**, 656–659.
- [11] Severcan, I., Geary, C., Verzemnieks, E., Chworos, A., and Jaeger, L. (Mar, 2009) Square-Shaped RNA Particles from Different RNA Folds. *Nanotechnology Letters*, **9**(3), 1270–1277.
- [12] Robertus, J. D., Ladner, J. E., Finch, J. T., Rhodes, D., Brown, R. S., Clark, B. F. C., and Klug, A. (1974) Structure of Yeast Phenylalanine tRNA at 3 Å Resolution. *Nature*, **250**, 546.
- [13] Kim, S. H. (1974) Three-Dimensional Tertiary Structure of Yeast Phenylalanine Transfer RNA. *Science*, **185**, 435.
- [14] Stout, C. D., Mizuno, H., Rubin, J., Brennan, T., Rao, S. T., and Sundaralingam, M. (Apr, 1976) Atomic Coordinates and Molecular Conformation of Yeast Phenylalanyl tRNA. An Independent Investigation. *Nucleic Acids Research*, **3**(4), 1111–1123.
- [15] Noller, H. F. (2005) RNA Structure: Reading the Ribosome. *Science*, **309**, 1508–1514.
- [16] Kruger, K., Grabowski, P. J., Zaug, A. J., Sands, J., Gottschling, D. E., and Cech, T. R. (1982) Self-Splicing RNA: Autoexcision and Autocyclization of the Ribosomal RNA Intervening Sequence of Tetrahymena. *Cell*, **31**, 147–157.



- [17] Guerrier-Takada, C., Gardiner, K., Marsh, T., Pace, N., and Altman, S. (1983) The RNA Moiety of Ribonuclease P is the Catalytic Subunit of the Enzyme. *Cell*, **35**, 849–857.
- [18] Zuker, M. (1989) On Finding All Suboptimal Foldings of an RNA Molecule. *Science*, **244**, 48–52.
- [19] Hofacker, I. L., Fontana, W., Stadler, P. F., Bonhoeffer, L. S., Tacker, M., and Schuster, P. (1994) Fast Folding and Comparison of RNA Secondary Structures. *Monatshefte fur Chemie*, **125**, 167–188.
- [20] Borer, P. N., Dengler, B., Tinoco, J. I., and Uhlenbeck, O. C. (1974) Stability of Ribonucleic Acid Double-Stranded Helices. *Journal of Molecular Biology*, **86**, 843–853.
- [21] Batey, R. T., Rambo, R. P., and Doudna, J. A. (1999) Tertiary Motifs in RNA Structure and Folding. *Angewandte Chemie International Edition*, **38**(16), 2326–2343.
- [22] Thirumalai, D. and Hyeon, C. (2005) RNA and Protein Folding: Common Themes and Variations. *Biochemistry*, **44**, 4957–4970.
- [23] Chen, S.-J. and Dill, K. A. (1995) Statistical Thermodynamics of Double-Stranded Polymer Molecules. *Journal of Chemical Physics*, **103**, 5802–5813.
- [24] Chen, S.-J. and Dill, K. A. (1998) Theory for the Conformational Changes of Double-Stranded Chain Molecules. *Journal of Chemical Physics*, **109**, 4602–4616.
- [25] Thirumalai, D. and Woodson, S. A. (1996) Kinetics of Folding of Proteins and RNA. *Accounts in Chemical Research*, **29**, 433–439.
- [26] Tinoco, I. and Bustamante, C. (1999) How RNA folds. *Journal of Molecular Biology*, **293**(2), 271–281.
- [27] Rangan, P., Masquida, B., Westhof, E., and Woodson, S. A. (2003) Assembly of Core Helices and Rapid Tertiary Folding of a Small Bacterial Group I Ribozyme. *Proceedings of the National Academy of Sciences of the United States of America*, **100**, 1574–1579.
- [28] Moore, P. B. The RNA World chapter The RNA Folding Problem, pp. 381–401 Cold Spring Harbor Laboratory Press 2nd edition (1999).
- [29] Sorin, E. J., Nakatani, B. J., Rhee, Y. M., Jayachandran, G., Vishal, V., and Pande, V. S. (2004) Does Native State Topology Determine the RNA Folding Mechanism?. *Journal of Molecular Biology*, **337**, 789–797.
- [30] Klein, D. J., Moore, P. B., and Steitz, T. A. (2004) The Contribution of Metal Ions to the Structural Stability of the Large Ribosomal Subunit. *RNA*, **10**(9), 1366–1379.
- [31] Malhotra, A., Tan, R. K., and Harvey, S. C. (1990) Prediction of the Three-Dimensional Structure of Escherichia Coli 30S Ribosomal Subunit: A Molecular Mechanics Approach. *Proceedings of the National Academy of Sciences of the United States of America*, **87**, 1950–1954.
- [32] Stagg, S. M., Mears, J. A., and Harvey, S. C. (2003) A Structural Model for the Assembly of the 30 S Subunit of the Ribosome. *Journal of Molecular Biology*, **328**, 49–61.
- [33] Parisien, M. and Major, F. (2008) The MC-Fold and MC-Sym Pipeline Infers RNA Structure from Sequence Data. *Nature*, **452**, 51–55.
- [34] Westhead, D., Slidel, T., Flores, T., and Thornton, J. (1999) Protein Structural Topology: Automated Analysis and Diagrammatic Representation. *Protein Science*, **8**, 897–904.

- [35] Gerstein, M. and Thornton, J. M. (2003) Sequences and Topology. *Current Opinion in Structural Biology*, **13**, 341–343.
- [36] Meiler, J. and Baker, D. (2003) Coupled Prediction of Protein Secondary and Tertiary Structure. *Proceedings of the National Academy of Sciences of the United States of America*, **100**, 12105–12110.
- [37] Bloomfield, V. A., Crothers, D. M., and Jr., I. T. (2000) *Nucleic Acids: Structures, Properties and Functions*, University Science Books, .
- [38] Boots, J. L., Canny, M. D., Azimi, E., and Pardi, A. (Oct, 2008) Metal Ion Specificities for Folding and Cleavage Activity in the Schistosoma Hammerhead Ribozyme. *RNA*, **14**(10), 2212–2222.
- [39] Zhuang, X. and Rief, M. (2003) Single-Molecule Folding. *Current Opinion in Structural Biology*, **13**, 88–97.
- [40] Liphardt, J., Onoa, B., Smith, S., Jr., I. T., and Bustamante, C. (2001) Reversible Unfolding of Single RNA Molecules by Mechanical Force. *Science*, **292**, 733–737.
- [41] Onoa, B. and Jr., I. T. (2004) RNA Folding and Unfolding. *Current Opinion in Structural Biology*, **14**(3), 374–379.
- [42] Tinoco, I. (2004) Force as a Useful Variable in Reactions: Unfolding RNA. *Annual Review of Biophysics & Biomolecular Structure*, **33**, 363–385.
- [43] Hyeon, C. and Thirumalai, D. (2005) Mechanical Unfolding of RNA Hairpins. *Proceedings of the National Academy of Sciences of the United States of America*, **102**(19), 6789–6794.
- [44] Crooks, G. E. (1999) Entropy Production Fluctuation Theorem and the Nonequilibrium Work Relation for Free-Energy Differences. *Physical Review E*, **60**, 2721–2726.
- [45] Collin, D., F.Ritort, Jarzynski, C., Smith, S. B., Jr., I. T., and Bustamante, C. (2005) Verification of the Crooks Fluctuation Theorem and Recovery of RNA Folding Free Energies. *Nature*, **437**, 231–234.
- [46] Wang, Y., Rader, A. J., Bahar, I., and Jernigan, R. L. (2004) Global Ribosome Motions Revealed with Elastic Network Model. *Journal of Structural Biology*, **147**, 302–314.
- [47] Bahar, I. and Jernigan, R. L. (1998) Vibrational Dynamics of Transfer RNAs: Comparison of the Free and Synthetase-Bound Forms. *Journal of Molecular Biology*, **281**, 871–884.
- [48] Wang, Y. and Jernigan, R. L. (2005) Comparison of tRNA Motions in the Free and Ribosomal Bound Structures. *Biophysical Journal*, **89**, 3399–3409.
- [49] Tung, C.-S. and Sanbonmatsu, K. Y. (2004) Atomic Model of the Thermus thermophilus 70S Ribosome Developed in Silico. *Biophysical Journal*, **87**, 2714–2722.
- [50] Sanbonmatsu, K. Y., Simpson, J., and Tung, C.-S. (2005) Simulating Movement of tRNA into the Ribosome During Decoding. *Proceedings of the National Academy of Sciences of the United States of America*, **102**, 15854–15859.
- [51] Sponer, J., Leszczynski, J., and Hobza, P. (1996) Nature of Nucleic Acid-Base Stacking: Nonempirical Ab Initio and Empirical Potential Characterization of 10 Stacked Base Dimers. Comparison of Stacked and H-Bonded Base Pairs. *Journal of Physical Chemistry*, **100**, 5590–5596.

- [52] Sponer, J., Leszczynski, J., and Hobza, P. (1997) Thioguanine and Thiouracil: Hydrogen-Bonding and Stacking Properties. *Journal of Physical Chemistry A*, **101**, 9489–9495.
- [53] Sponer, J., Berger, I., Spackova, N., Leszczynski, J., and Hobza, P. (2000) Aromatic Base Stacking in DNA: From Ab Initio Calculations to Molecular Dynamics Simulations. *Journal of Biomolecular Structure and Dynamics*, **11**, 1–24.
- [54] Sponer, J., Jureka, P., Marchan, I., Luque, F. J., Orozco, M., and Hobza, P. (2006) Nature of Base Stacking: Reference Quantum-Chemical Stacking Energies in Ten Unique B-DNA Base-Pair Steps. *Chemistry - A European Journal*, **12**, 2854–2865.
- [55] Hobza, P. and Sponer, J. (2002) Toward True DNA Base-Stacking Energies: MP2, CCSD(T), and Complete Basis Set Calculations. *Journal of the American Chemical Society*, **124**, 11802–11808.
- [56] Marky, L. A. and Breslauer, K. J. (1982) Calorimetric Determination of Base-Stacking Enthalpies in Double-Helical DNA Molecules. *Biopolymers*, **11**, 2185–2194.
- [57] Manning, G. S. (1977) Limiting Laws and Counterion Condensation in Polyelectrolyte Solutions IV. The Approach to the Limit and the Extraordinary Stability of the Charge Fraction. *Biophysical Chemistry*, **7**, 95–102.
- [58] Manning, G. S. (2003) Comments on Selected Aspects of Nucleic Acid Electrostatics. *Biopolymers*, **69**, 137–143.
- [59] Antypov, D., Barbosa, M. C., and Holm, C. (2005) Incorporation of Excluded-Volume Correlations into Poisson-Boltzmann Theory. *Physical Review E*, **71**(6), 1–6.
- [60] Berman, H. M., Olson, W. K., Beveridge, D. L., Westbrook, J., Gelbin, A., Demeny, T., Hsieh, S. H., Srinivasan, A. R., and Schneider, B. (1992) The Nucleic Acid Database. A Comprehensive Relational Database of Three-Dimensional Structures of Nucleic Acids. *Biophysical Journal*, **63**, 751–759.
- [61] Murray, L. J. W., III, W. B. A., Richardson, D. C., and Richardson, J. S. (2003) RNA Backbone is Rotameric. *Proceedings of the National Academy of Sciences of the United States of America*, **100**, 13904–13909.
- [62] Schneider, B., Moravek, Z., and Berman, H. (2004) RNA Conformational Classes. *Nucleic Acids Research*, **32**, 1666–1677.
- [63] Sykes, M. T. and Levitt, M. (2005) Describing RNA Structure by Libraries of Clustered Nucleotide Doublets. *Journal of Molecular Biology*, **351**, 26–38.
- [64] Richardson, J. S., Schneider, B., Murray, L. W., Kapral, G. J., Immormino, R. M., Headd, J. J., Richardson, D. C., Ham, D., HersHKovits, E., Williams, L. D., Keating, K. S., Pyle, A. M., Micalef, D., Westbrook, J., and Berman, H. M. (2008) RNA Backbone: Consensus All-Angle Conformers and Modular String Nomenclature (An RNA Ontology Consortium Contribution). *RNA*, **14**, 465–481.
- [65] Zuker, M. (2003) Mfold Web Server for Nucleic Acid Folding and Hybridization Prediction. *Nucleic Acids Research*, **31**(13), 3406–3415.
- [66] Mathews, D. H. and Turner, D. H. (Mar, 2002) Dynalign: An Algorithm for Finding the Secondary Structure Common to Two RNA Sequences. *Journal of Molecular Biology*, **317**(2), 191–203.
- [67] Chen, S.-J. and Dill, K. A. (2000) RNA Folding Energy Landscapes. *Proceedings of the National Academy of Sciences of the United States of America*, **97**, 646–651.

- [68] Das, R. and Baker, D. (Sep, 2007) Automated de Novo Prediction of Native-Like RNA Tertiary Structures. *Proceedings of the National Academy of Sciences of the United States of America*, **104**(37), 14664–14669.
- [69] Ding, F., Sharma, S., Chalasani, P., Demidov, V. V., Broude, N. E., and Dokholyan, N. V. (Jun, 2008) Ab Initio RNA Folding by Discrete Molecular Dynamics: From Structure Prediction to Folding Mechanisms. *RNA*, **14**(6), 1164–1173.
- [70] Jonikas, M. A., Radmer, R. J., and Altman, R. B. (Dec, 2009) Knowledge-Based Instantiation of Full Atomic Detail Into Coarse-Grain RNA 3D Structural Models. *Bioinformatics*, **25**(24), 3259–3266.
- [71] Martinez, H. M., Jr, J. V. M., and Shapiro, B. A. (2008) RNA2D3D: A Program for Generating, Viewing, and Comparing 3-Dimensional Models of RNA. *Journal of Biomolecular Structure and Dynamics*, **25**, 573–752.
- [72] Richardson, J. S. (2000) Early Ribbon Drawings of Proteins. *Nature Structural Biology*, **7**, 624–625.
- [73] Richardson, D. C. and Richardson, J. S. (2002) Teaching Molecular 3-D Literacy. *Biochemistry and Molecular Biology Education*, **30**, 21–26.
- [74] Andreeva, A., Howorth, D., Brenner, S. E., Hubbard, T. J. P., Chothia, C., and Murzin, A. G. (2004) SCOP Database in 2004: Refinements Integrate Structure and Sequence Family Data. *Nucleic Acids Research*, **32**, D226–D229.
- [75] Klosterman, P. S., Tamura, M., Holbrook, S. R., and Brenner, S. E. (2002) SCOR: a Structural Classification of RNA Database. *Nucleic Acids Research*, **30**, 392–394.
- [76] Klosterman, P. S., Hendrix, D. K., Tamura, M., Holbrook, S. R., and Brenner, S. E. (2004) Three-Dimensional Motifs from the SCOR, Structural Classification of RNA Database: Extruded Strands, Base Triples, Tetraloops and U-turns. *Nucleic Acids Research*, **32**(8), 2342–2352.
- [77] Holbrook, S. R. (2005) RNA Structure: The Long and the Short of it. *Current Opinion in Structural Biology*, **15**, 302–308.
- [78] Leontis, N. B. and Westhof, E. (2003) Analysis of RNA Motifs. *Current Opinion in Structural Biology*, **13**, 300–308.
- [79] Moore, P. B. (1999) Structural Motifs in RNA. *Annual Review of Biochemistry*, **68**, 287–300.
- [80] Nasalean, L., Stombaugh, J., Zirbel, C. L., and Leontis, N. B. Vol. 13, of Springer Series in Biophysics chapter Chapter I, pp. 1–26 Springer Verlag Berlin Heidelberg (November, 2009).
- [81] Lemieux, S. and Major, F. (2006) Automated Extraction and Classification of RNA Tertiary Structure Cyclic Motifs. *Nucleic Acids Research*, **34**, 2340–2346.
- [82] Duarte, C. M., Wadley, L. M., and Pyle, A. M. (2003) RNA Structure Comparison, Motif Search and Discovery Using a Reduced Representation of RNA Conformational Space. *Nucleic Acids Research*, **31**(16), 4755–4761.
- [83] Laing, C., Jung, S., Iqbal, A., and Schlick, T. (Oct, 2009) Tertiary Motifs Revealed in Analyses of Higher-Order RNA Junctions. *Journal of Molecular Biology*, **393**(1), 67–82.
- [84] Laing, C. and Schlick, T. (Jul, 2009) Analysis of Four-way Junctions in RNA Structures. *Journal of Molecular Biology*, **390**(3), 547–559.

- [85] Holbrook, S. R. (2008) Structural Principles From large RNAs. *Annual Review in Biophysics*, **37**, 445–464.
- [86] Spacková, N. and Sponer, J. (2006) Molecular Dynamics Simulations of Sarcin-Ricin rRNA Motif. *Nucleic Acids Research*, **34**(2), 697–708.
- [87] Réblová, K., Spacková, N., Stefl, R., Csaszar, K., Koca, J., Leontis, N. B., and Sponer, J. (Jun, 2003) Non-Watson-Crick Basepairing and Hydration in RNA Motifs: Molecular Dynamics of 5S rRNA Loop E. *Biophysical Journal*, **84**(6), 3564–3582.
- [88] Sarver, M., Zirbel, C. L., Stombaugh, J., Mokdad, A., and Leontis, N. B. (Jan, 2008) FR3D: Finding Local and Composite Recurrent Structural Motifs in RNA 3D Structures. *Journal of Mathematical Biology*, **56**(1-2), 215–252.
- [89] Mokdad, A. and Frankel, A. D. (April, 2008) ISFOLD: Structure Prediction of Base-pairs in Non-helical RNA Motifs From Isostericity Signatures in Their Sequence Alignments. *Journal of Biomolecular Structure and Dynamics*, **25**(5), 467–472.
- [90] Stonge, K., Thibault, P., Hamel, S., and Major, F. (2007) Modeling RNA Tertiary Structure Motifs by Graph-Grammars. *Nucleic Acids Research*, 2007, 11, pp. 1–11.

## Chapter 2

### RNA Base Steps

The problem of classification of the space of conformations of RNA is not new, see for example, Olson 1972 [1], Saenger 1984 [2], and Gautheret 1993 [3]. This problem had only been addressed by a few researchers before the turn of the twenty first century, but starting in the year 2000 a vast amount of RNA structural information has become available with the elucidation of the structure of the 30S small ribosomal subunit of *Thermus thermophilus*, a bacterial ribosome [4, 5], and the 50S large ribosomal subunit of *Haloarcula marismortui*, an archaeal ribosome [6].

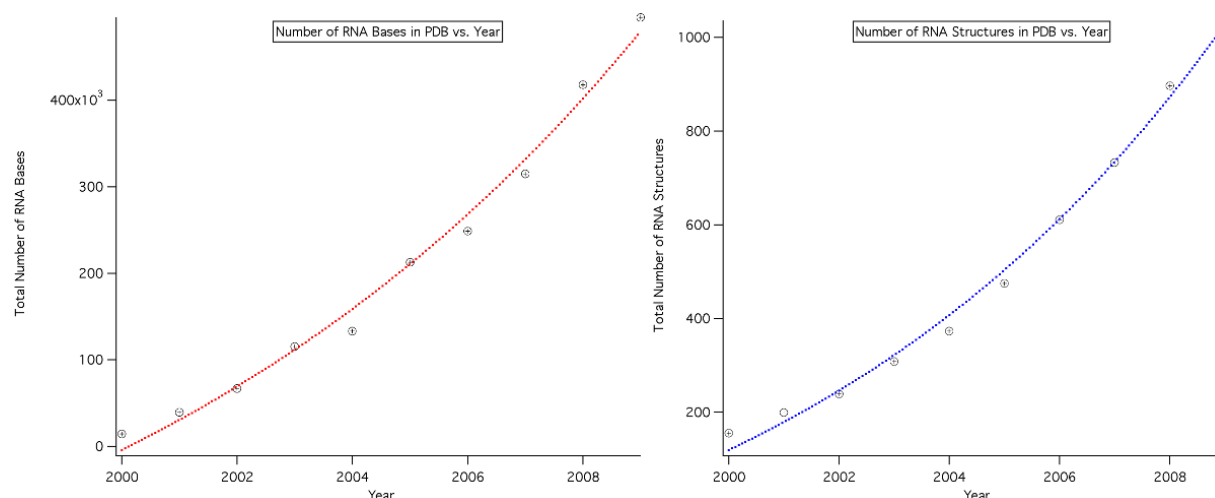


Figure 2.1: **Right:** Total number of RNA bases added to the PDB database between 2000 and 2010 (Exponential fit line in blue). **Left:** Total number of RNA structures solved yearly by X-Ray crystallography between 2000 and 2010 (Exponential fit line in red).

Between 1978 and 2000 a total of 116 RNA structures with resolution greater than 3.5Å, and comprising around 5500 nucleotide bases are found in the Protein Data Bank (PDB), and between 2000 and today a total of 931 RNA structures comprising 491158 nucleotide bases are found. That is, the increase in information due to the solution of large RNA structures is about two orders of magnitude as pointed out by Noller [7]. Looking at the growth of RNA structural information from 2000 until today, it is clear that both the total number of RNA structures deposited to the PDB, and the total number of nucleotide bases in these structures, is growing in an exponential way (as can be seen by the exponential fits in Figure 2.1). It's important to note that such growth comes mainly from ribosomal structures which contain 88 percent of all RNA bases in the PDB. So, even though structural interest in RNA is growing since ribosomal structures became available in 2000, and several Nobel prizes have been awarded for work in this field, along with the exciting possibilities of deciphering large RNA [8] structures other than the ribosome, still the growth of the RNA structural field is far from that of proteins if weighed by the growth in diversity of RNA structural information in the past decade. At the present time if we look at the distribution of RNA sizes counted by number of bases, as can be seen in Figure 2.2 it's clear that there are great patches where there are no RNA structures whatsoever, roughly between 600 and 1400 bases and between 1800 and 2700 bases. The area of non-coding RNA's holds great promise

for finding structured RNA's in such length ranges as has recently been suggested by Breaker [8] A representative example of the characteristic ranges of RNA structures available to date in the PDB can be seen in Table 2 for structures larger than 300 bases.

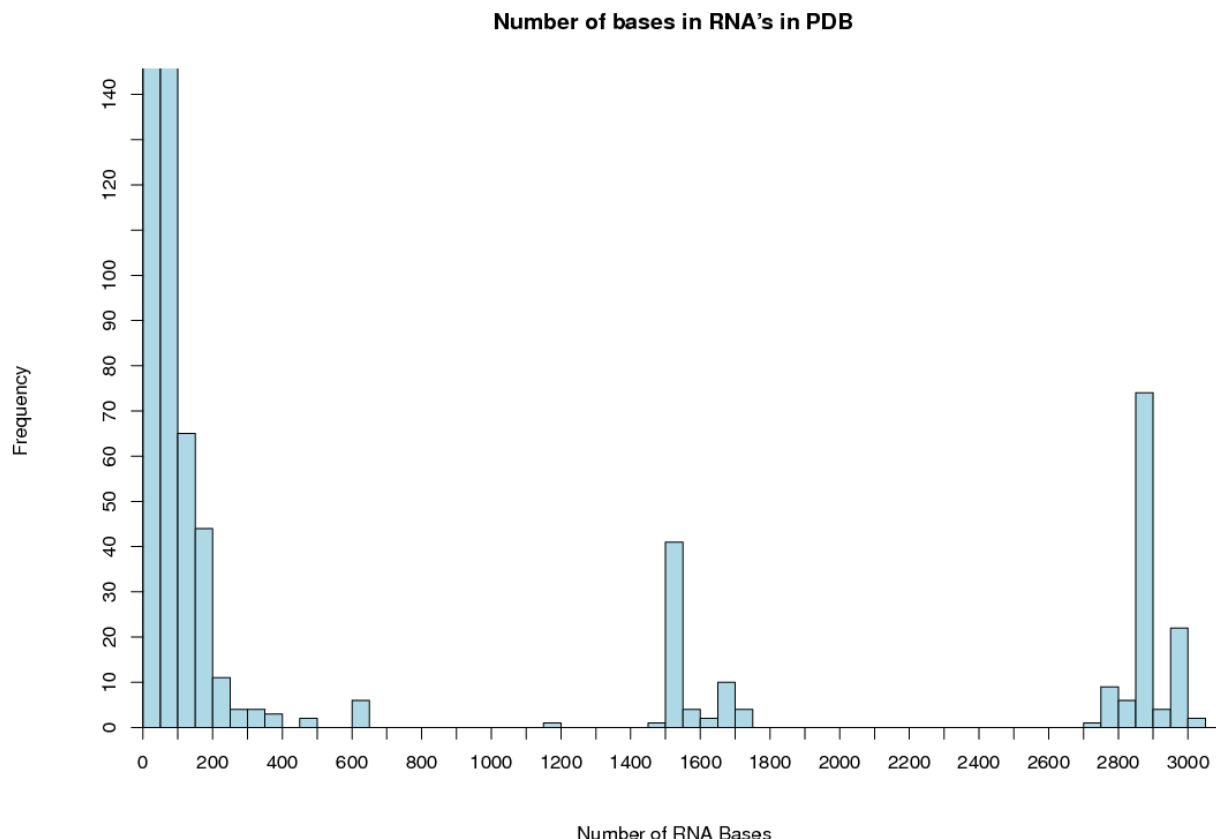


Figure 2.2: Frequency of nucleotide bases in RNA molecules found in the PDB classified by the size of RNA molecules. We define the size as the total number of nucleotide bases present per molecule.

The analysis of RNA conformational information contained in RNA structural data can be divided into three main perspectives: an atom based perspective; a bond based perspective; and a third, as yet unexplored to our knowledge, rigid-body based perspective. In the atom based perspective, either direct comparison of backbone atom positions is made [9], or a comparison of distances between a reduced set of atoms taken from the nucleotide backbone, sugar, and base [10]. The bond based perspective is divided into three main categories; the first considers the consecutive covalent bonds in the RNA backbone and the glycosidic bond between the sugar and base, that is, six backbone torsion angles and one glycosidic torsion angle [9, 11, 12, 13, 14]; or alternatively the pseudo-bonds between consecutive P and C4' atoms and the resulting pseudo-torsion angles  $\eta$  and  $\theta$  [1, 15, 16, 17]. The third category considers the networks of horizontal hydrogen bonding patterns coming from a definition of interacting edge boundaries in the nucleotide bases [18, 19, 20]. In this chapter we study the rigid body based perspective using clustering analysis.

## 2.1 Consensus Clustering of Single Stranded Base Step Parameters

To our knowledge there has been no classification of rigid-body base-step parameters for RNA structures deposited at the PDB. It is important to note here that in crystal structures, RNA bases are

PDBID	Structure Name	Phylogenetic Group	Number of bases	Year
1l8v	Mutant of P4-P6 Domain of Group I Intron	Eukaryote	314	2002
3igi	Group II Intron	Bacteria	395	2009
1fg0	Central Loop in Domain V of 23S rRNA	Archaea	499	2000
2nz4	GlmS Ribozyme	Eukaryote	604	2006
1xmq	30S rRNA	Bacteria	1522	2004
1ffk	50S rRNA Subunit	Archaea	2828	2000

Table 2.1: Some large RNA structures (>300 bases) elucidated in the last decade.

determined more accurately than backbone torsion angles, as has been shown by Richardson and collaborators from analysis of van der Waals steric clashes. This can be seen more clearly in Figure 2.3, reproduced from Richardson's work [11], where the red and orange dots in the backbone atoms region denote steric clashes and the green and yellow dots in the base atoms region denote very good agreement with expected van der Waals distances.

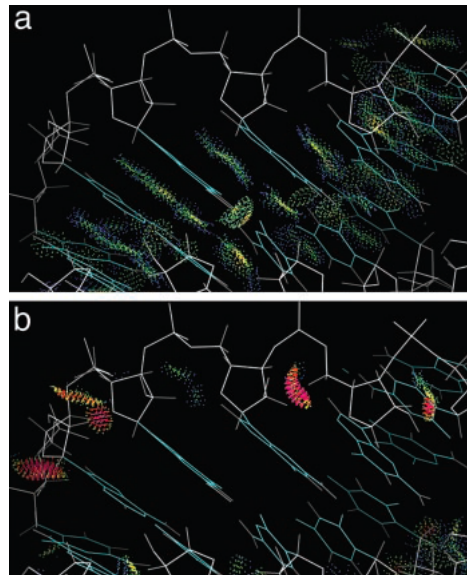


Figure 2.3: Figure taken from Richardson et al. [11] where the blue and green dots in a) mean very accurate van der Waals distances, and in b) the red and orange dots mean steric clashes, that is, distances outside the acceptable van der Waals range.

### 2.1.1 Combining Fourier Averaging Results and Clustering Analysis

Using the coordinates files of 20 rRNA structures provided by Schneider et al. [13] we used standard clustering analysis (CA) techniques (see Appendix A) to classify a set of non-ARNA base-steps using, rather than the more common torsion angles space, the base-step parameters space, that is, three translational parameters (Shift  $D_x$ , Slide  $D_y$ , Rise  $D_z$ ), and three rotational parameters (Tilt  $\tau$ , Roll  $\rho$ , Twist  $\omega$ ), which we describe with the hexaparametric vector  $\nu$ :

$$\nu = (D_x, D_y, D_z, \tau, \rho, \omega) \quad (2.1)$$



The results illustrated in Figures 2.4 C.1 and 2.5 were obtained by performing clustering analysis and consensus clustering on 20 structures provided by Schneider et al. [13]. These twenty structures were obtained by Schneider applying a Fourier averaging technique, and lexicographical clustering, to torsion angles of 23S rRNA. The methodology we used follows that used by others to recover the periodic table classification from multidimensional property vectors for elements [21, 22].

Group I contains a single structure 1 with base-plane normals pointing in opposite directions, Group II includes extended conformations with neighboring bases roughly parallel but not stacked and is formed by structures 15, 16, 10, 14, Group III also contains extended conformations with bases perpendicular to one another and is formed by structures 8, 9, 17, Group IV 18, 19, 20, 13, 11, 12, 5, 3, 6, 7, 2, 4 contains four major subgroups: (a) structures 2, 4 which are unstacked with bases neither parallel nor perpendicular; (b) structures 18, 19, 20 which are A-RNA related; (c) structures 11, 12, 13 which are unstacked and have parallel bases; and (d) structures 3, 5, 6, 7 which are also unstacked and have parallel bases. We also see in Group IV that the conformers in subgroups IV (c) and IV (d) are closely related and that the dimers in these two subgroups are more closely related to those in subgroup IV (b) than to those in subgroup IV (a).

When looking at Table 2.2, it's clear that there are 1858 steps (67%) which are not classified into any of the groups. The reason for this is the mixing of Fourier averaging for backbones, and the base step perspective. It might also be that we are not using the other A-RNA like backbone based structures from schneider's paper.

Right now I am doing a validation with clValid, to see if anything pops up regarding the "optimal" number of clusters for the data. Perhaps it would be wise to filter the data by proximity to A-RNA like conformation, say, take all structures which are some RMSD, or manhattan, or euclidean distance apart from the canonical A-RNA step parameters which are in table such and such.

Leave this argument for second part.

Table 2.2 shows the residue numbers of bases from 23S rRNA which belong to the main categories of Figure 2.5. To match residues of 23S rRNA belonging to the non-Atype clusters, a root mean squared deviation (RMSD) of 15 or less was required between step parameter vectors of 23S rRNA and the mean parameter vectors for the four non-Atype groups identified.

## 2.1.2 Partitional Clustering for Rigid Body Parameters

The argument I thought could have been made was that with clustering analysis alone on the whole data set, the A-RNA data would split naturally, without recurring to other ideas like Fourier Filtering of Bergman et al.

We also analyze the 2753 base-step parameter vectors in the ribosome. For the partitional clustering case, again, there is no known number of clusters in which the data must group, therefore we've calculated the within clusters sum of squares and also the average silhouette widths, for a particular selection of the number of partitions of the data for  $k = [2 - 80]$ . From figure 2.6 we can't conclude much. We see that the value of the within clusters sum of squares becomes constant around  $k = 47$  and there's also a change of curvature around  $k = 13$ . For the case where the average silhouette width has been computed, that is, figure 2.7, we see that the maximum is for  $k = 2$ , and there are some interesting maxima at  $k = 9, 12$ . Now that we have a clue as to which number of partitions the data optimally has we have plotted the k-means results for  $k = 13$  and  $k = 47$  in Figures number and number, and the PAM results for  $k = 2, 9, 12$  in Figure number.

We have also filtered the data according to the 16 possible RNA base steps, that is, AA, AG, GA, GG, UU, UC, CU, CC, UA, UG, CA, CG, AU, AC, GU, and GC. Tables showing how many representatives steps there are belonging to non-helical, helical, and watson-crick sets, will be later included and discussed here.

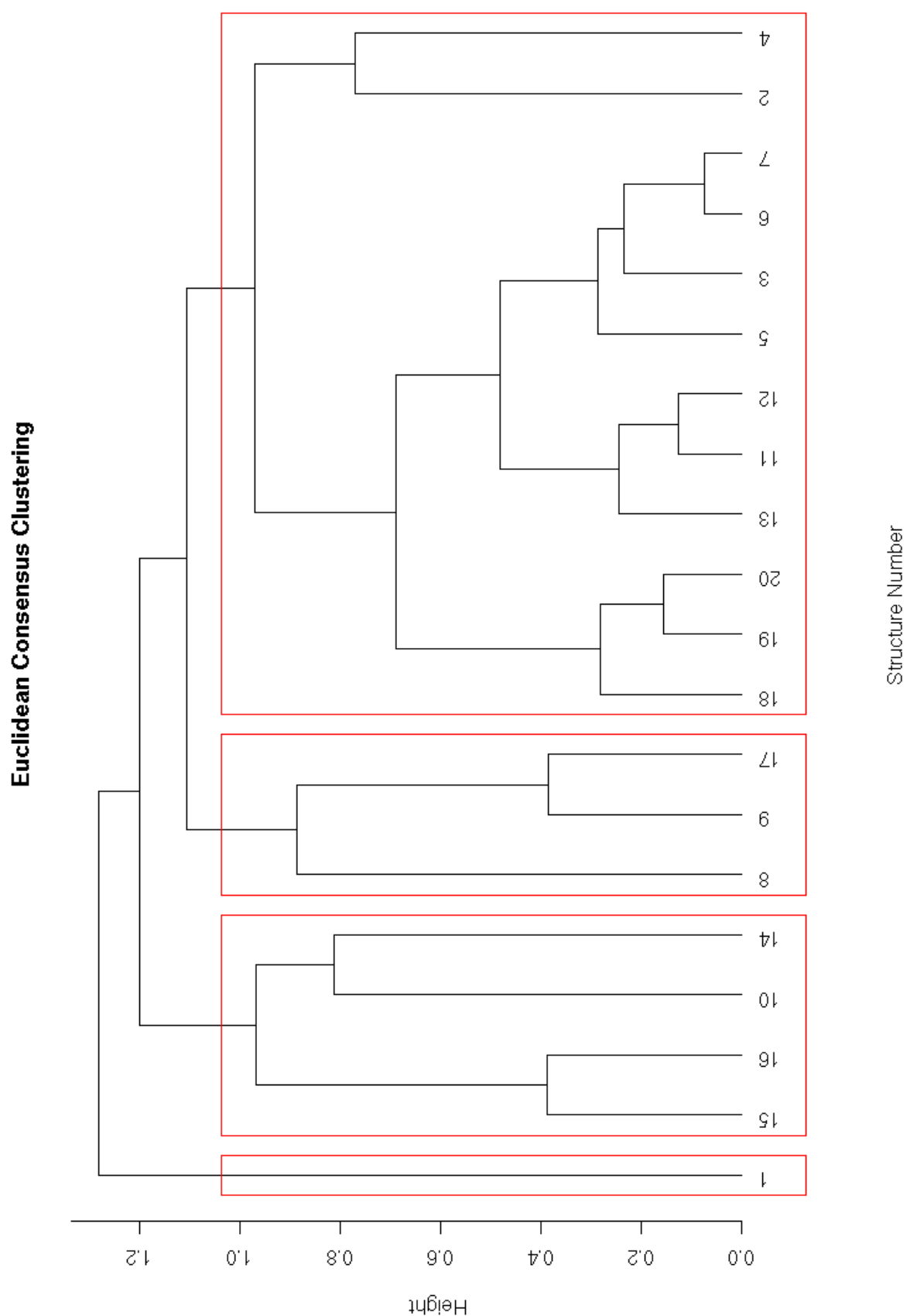


Figure 2.4: Dendrogram showing the results of consensus clustering of 20 non-Atype rRNA dinucleotides according to their hexadimensional base-step parameter vectors.

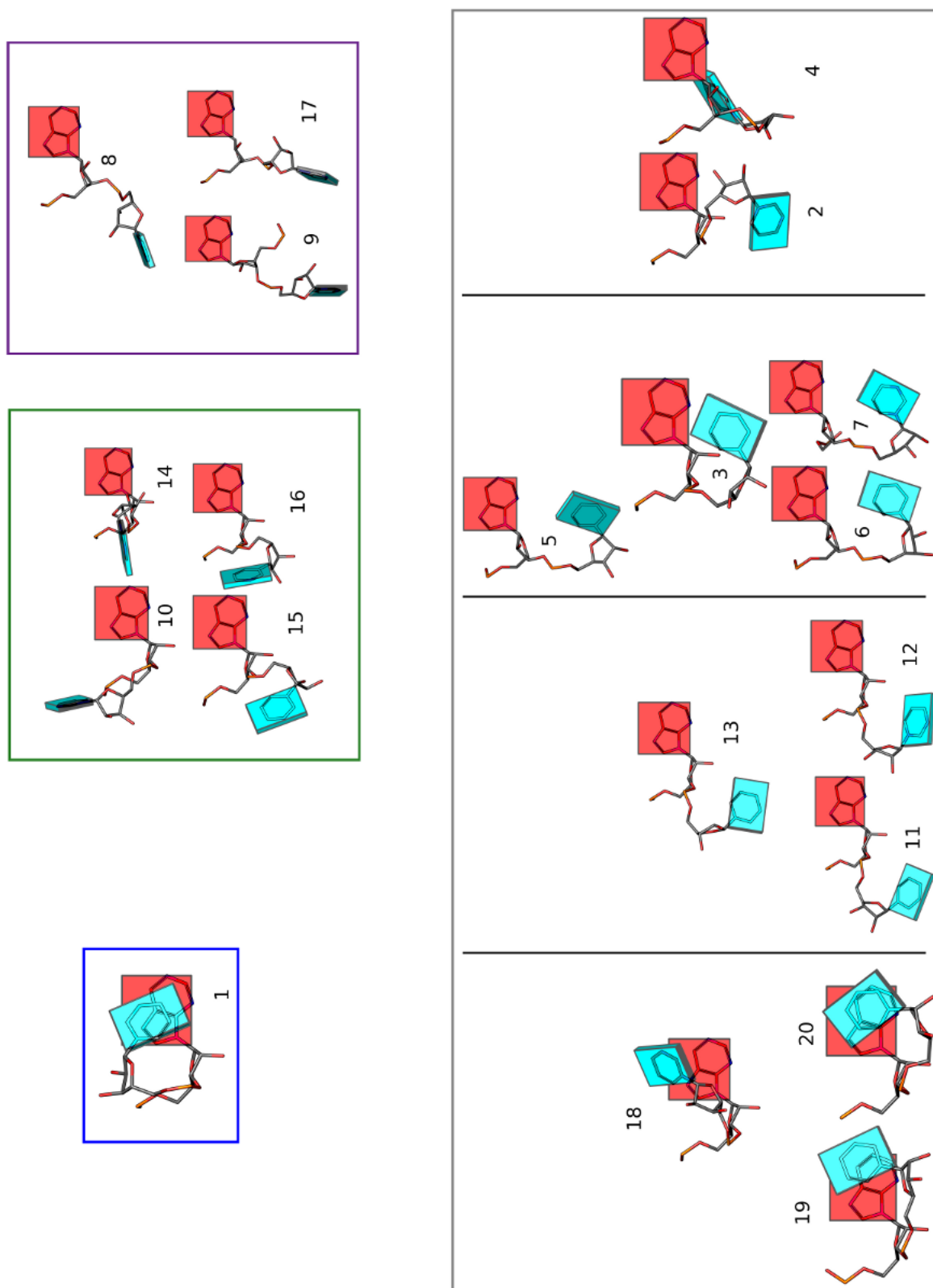


Figure 2.5: rRNA dinucleotide structures organized by clusters obtained from consensus clustering of their hexadimensional base-step parameter vectors.

Total Number of Nucleotides	RMSD Limit	Group	Base-steps	Base-step Residue Number	Overlaps
2754	< 15	I	3	892, 2006, 2390	
		II	5	459, 1279, 1653, 1919, 2302	
		III	1	2109	
		IV	35	79, 112, 128, 190, 213, 269, 358, 434, 488, 564, 706, 720, 775, 867, 966, 1292, 1503, 1543, 1614, 1766, 1874, 1908, 1971, 2017, 2257, 2427, 2516, 2540, 2755, 2782, 2810, 2826, 2874, 2882, 2913	
		IVa	1	882	
		IVb	807		
		IVc	9	306, 789, 854, 880, 1107, 1192, 1493, 1818, 2005	
		IVd	35	175, 213, 246, 264, 304, 358, 464, 518, 531, 534, 588, 795, 938, 1214, 1231, 1316, 1340, 1370, 1605, 1745, 1766, 1971, 1976, 2010, 2017, 2291, 2320, 2428, 2469, 2481, 2516, 2532, 2755, 2826, 2882	Only IVd with IV (213, 358, 1766, 1971, 2017, 2516, 2755, 2826, 2882)

Table 2.2: Residue numbers for base-steps with RMSD values less than 15 between the reference base-step vectors from the four groups of non-A-type RNA dinucleotide conformations and all base-step vectors found in the 23S strand of *Haloarcula marismortui* large ribosomal subunit.

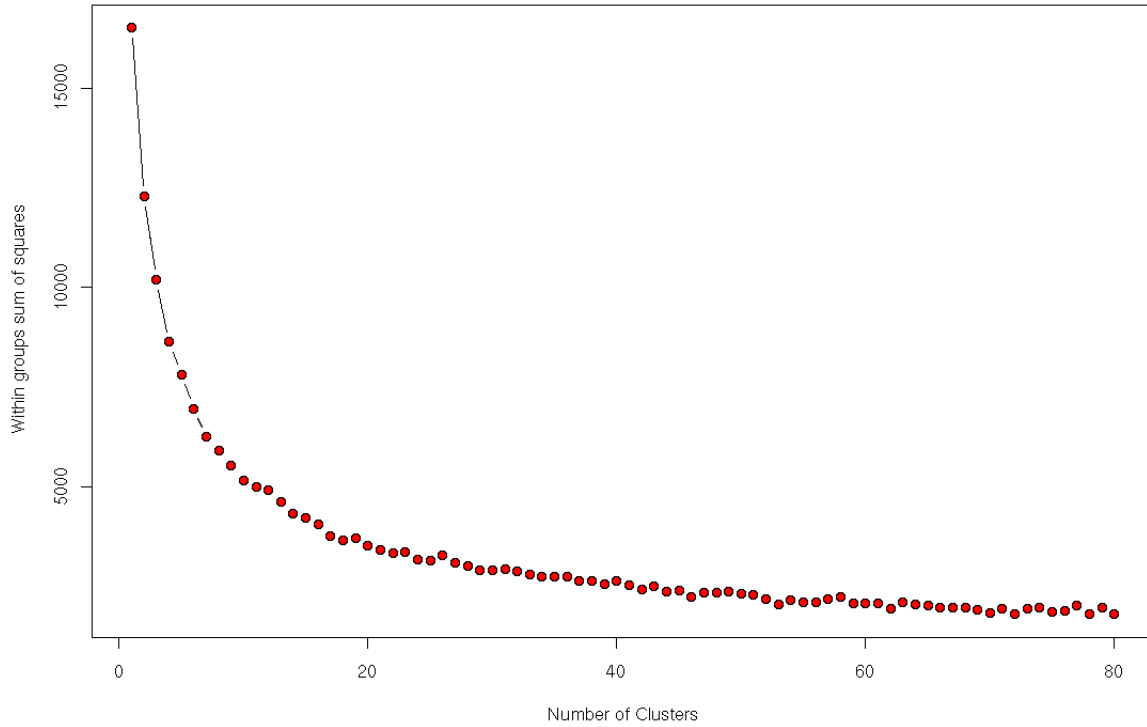


Figure 2.6: Sum of all within clusters sum of squares against number of clusters.

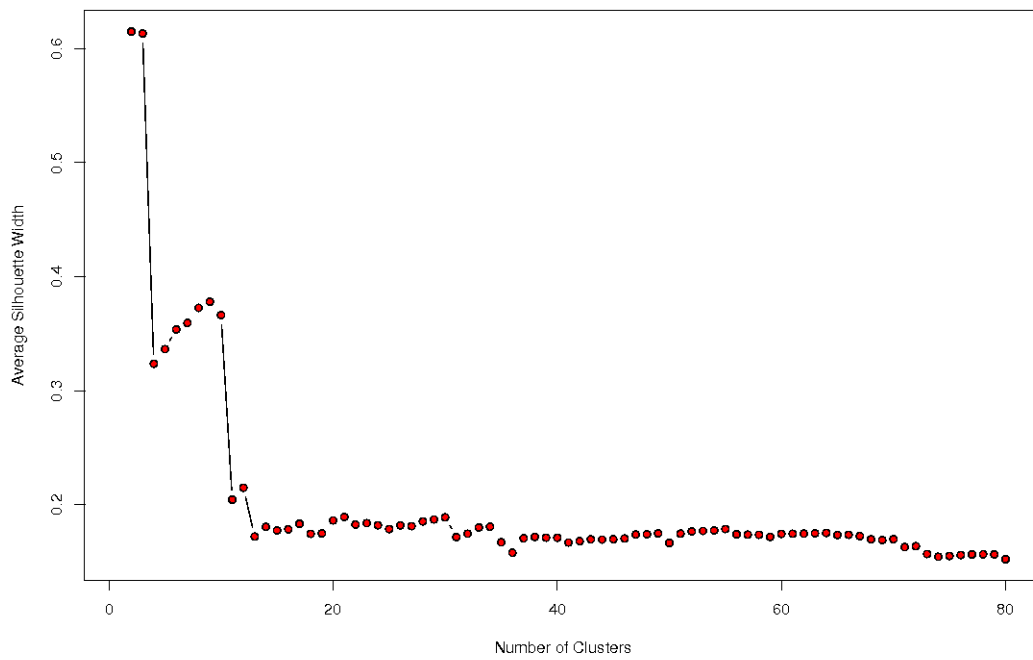


Figure 2.7: Average silhouette width against number of clusters.

### 2.1.3 Hierarchical Clustering for Rigid Body Parameters

Also as has been carried out for torsion angles, hierarchical clustering has also been performed on rigid body parameters, the results are yet to be included here. A cluster dissimilarity tree can be seen in Figure 2.8 for the 12 trees resulting from the four clustering methods and three distance definitions used to cluster the base step data.

## 2.2 RNA Conformations

There are two main RNA conformations, A-RNA ,and A'RNA, and maybe even a third unconfirmed one A"RNA [2]. Their values for their standard torsion angles and step parameters can be seen in Tables 2.3 and 2.4

Structure Name	$\alpha$	$\beta$	$\gamma$	$\delta$	$\epsilon$	$\zeta$	$\chi$	Reference
A-RNA	-68.9	179.5	54.5	82.2	-153.9	-70.8	-161.1	Arnott
A'-RNA	-70.0	176.6	60.8	76.7	-153.4	-69.4	-163.4	Arnott
AII-RNA	-65.0	175.1	52.9	81.1	-166.0	-68.0	-157.0	Schneider

Table 2.3: Base step torsion angles for the different known RNA conformations.

Structure Name	Shift ( $D_x$ )	Slide ( $D_y$ )	Rise ( $D_z$ )	Tilt ( $\tau$ )	Roll ( $\rho$ )	Twist ( $\Omega$ )	Reference
A-DNA	0.36	-1.39	3.29	2.46	12.50	30.19	
B-DNA	0.44	0.47	3.33	4.63	1.77	35.67	
A-RNA	-0.08	-1.48	3.30	-0.43	8.64	31.57	Arnott
A'-RNA	0.05	-1.88	3.39	-0.12	5.43	29.52	Arnott
AII-RNA	1.01	-2.52	3.33	2.94	9.75	25.12	Schneider

Table 2.4: Base step parameters for the different known RNA conformations. Notice that the base step parameters are for single bases rather than base-pairs.

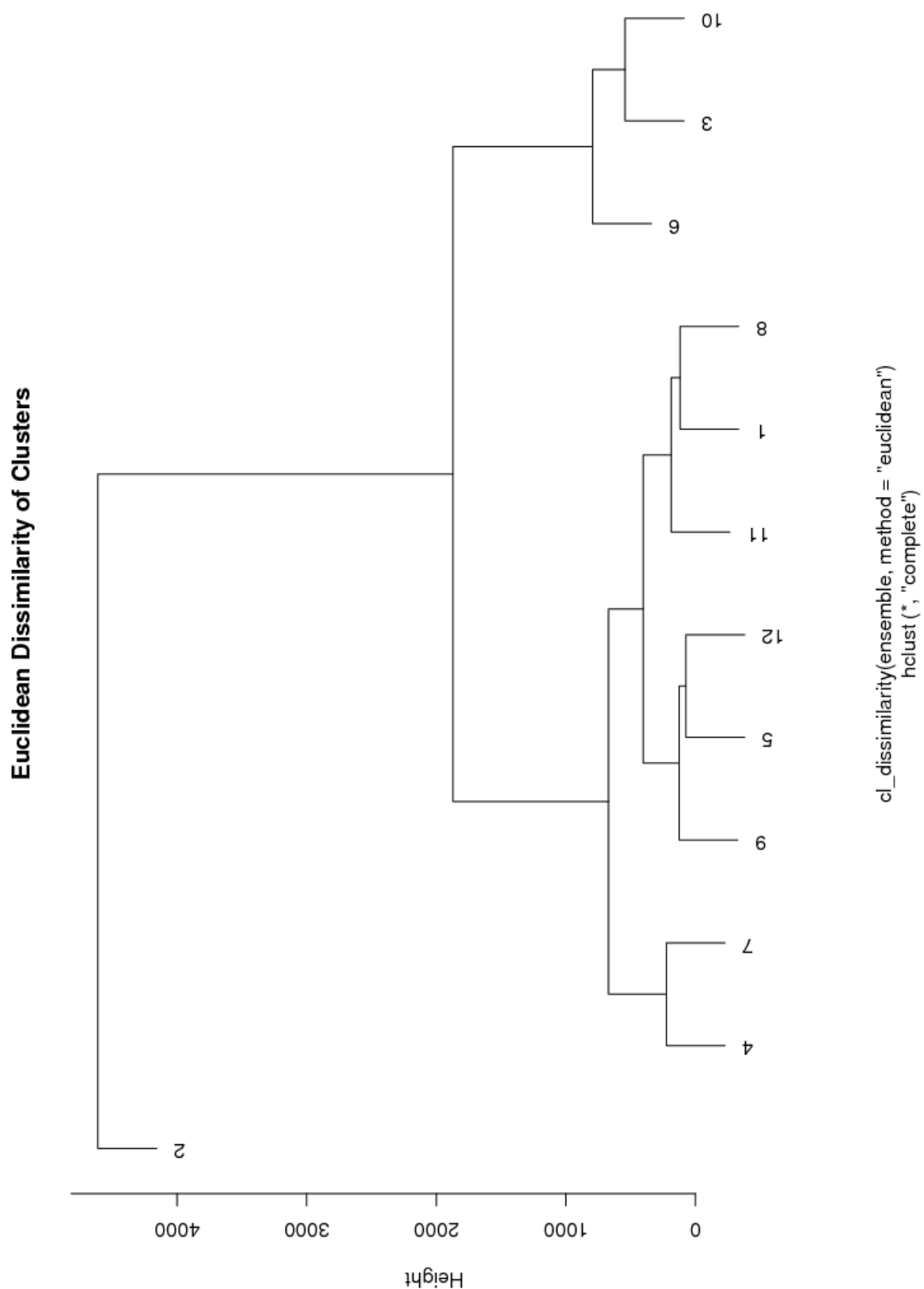


Figure 2.8: Cluster dissimilarities for the twelve hierarchical trees obtained from clustering of the six-dimensional base-step parameters obtained from the large subunit of the ribosome (PDB-ID:1jj2)

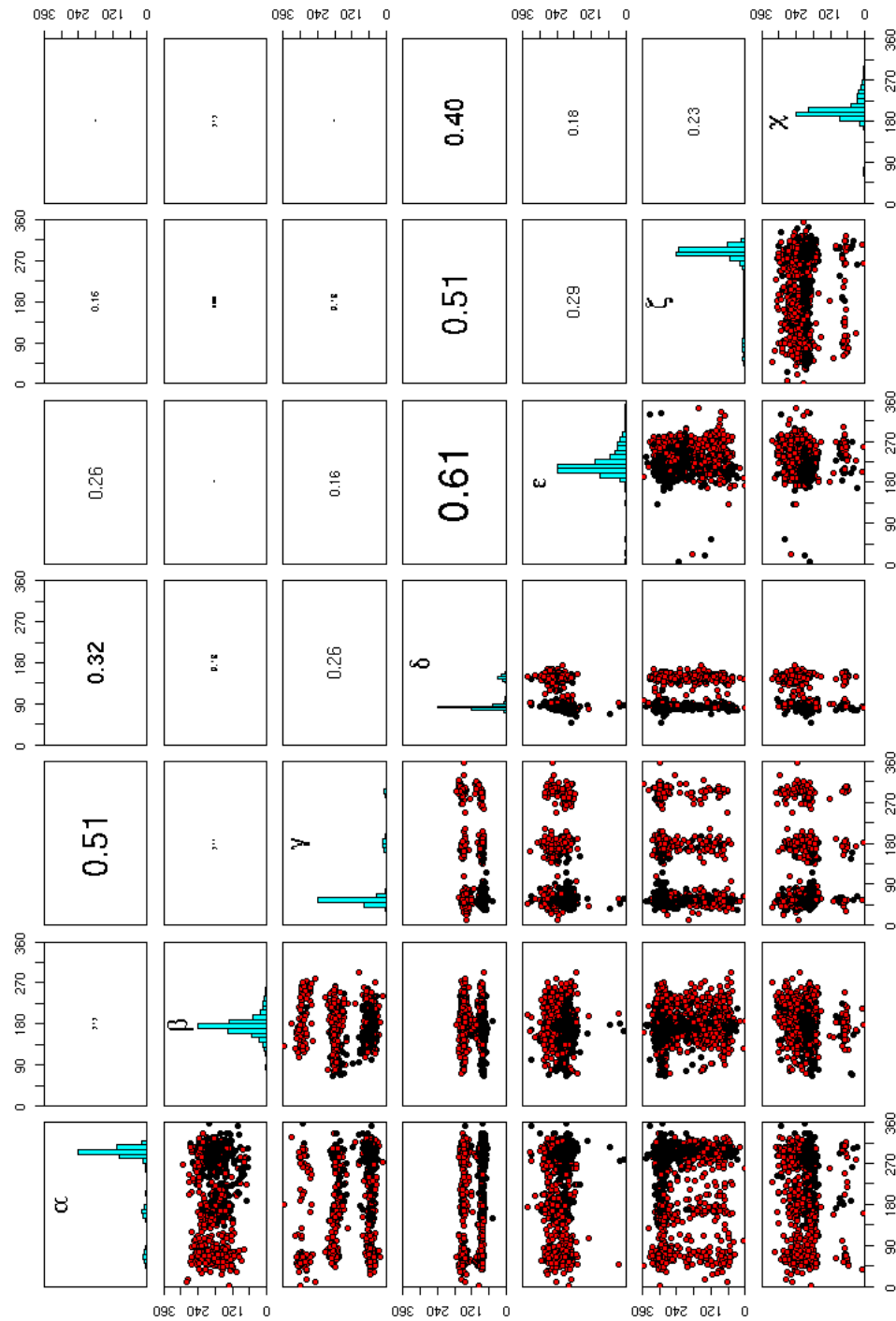


Figure 2.9: K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the *Hartigan-Wong* algorithm. The number of partitions is 2. The upper diagonal matrix displays the values of the linear correlation coefficient  $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal.



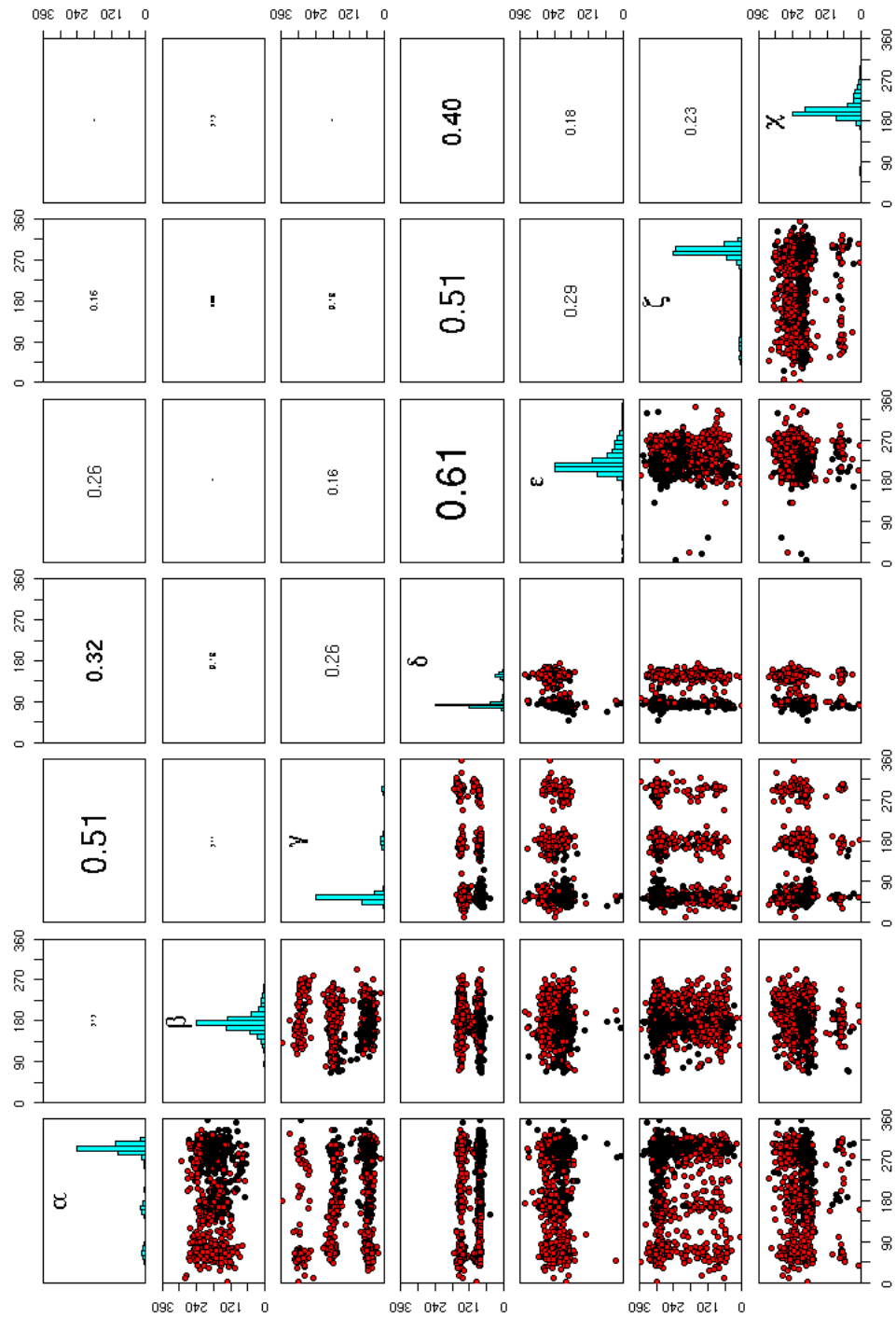


Figure 2.10: K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the *Lloyd* algorithm. The number of partitions is 2. The upper diagonal matrix displays the values of the linear correlation coefficient  $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal.

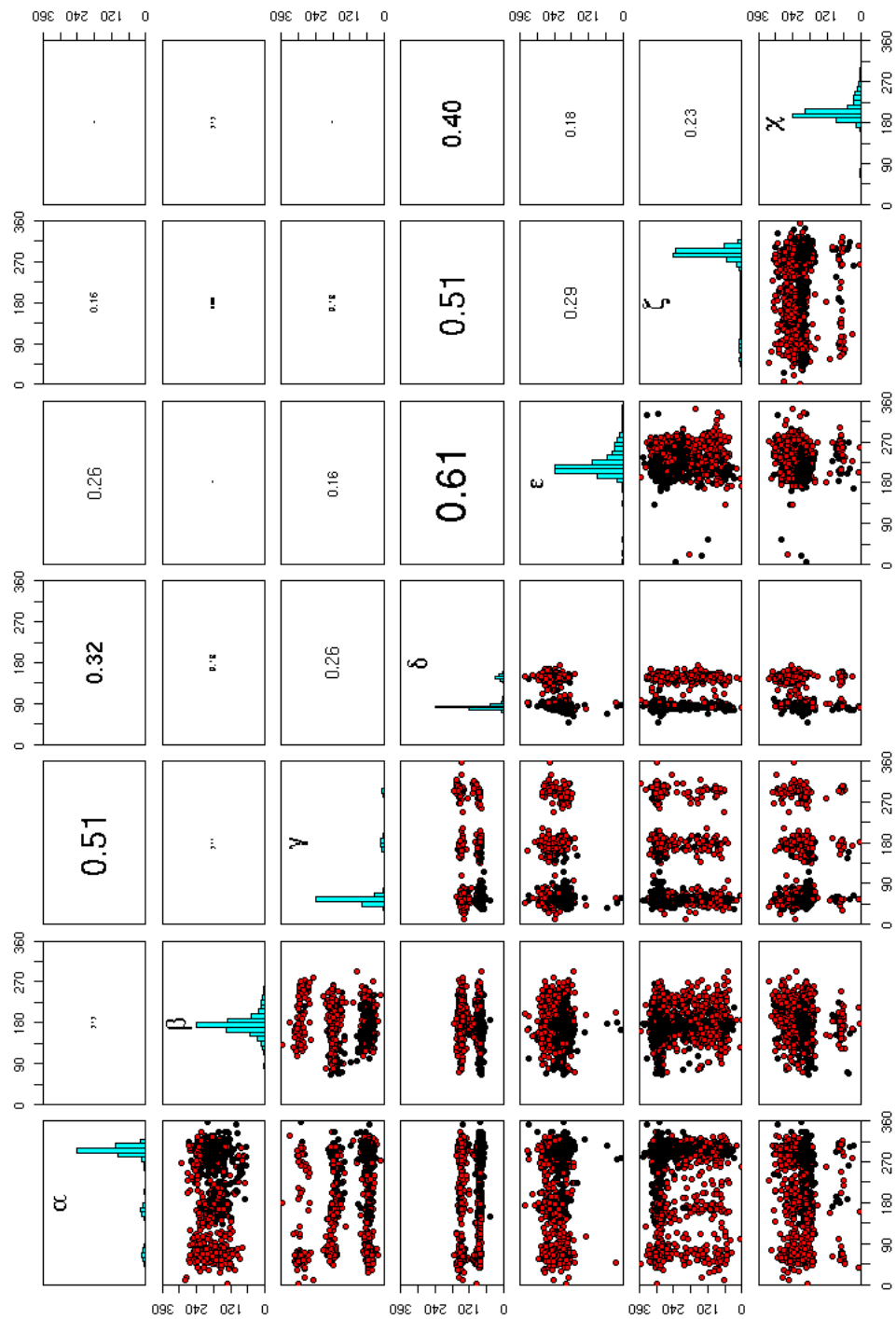


Figure 2.11: K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the *Forgy* algorithm. The number of partitions is 2. The upper diagonal matrix displays the values of the linear correlation coefficient  $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal.

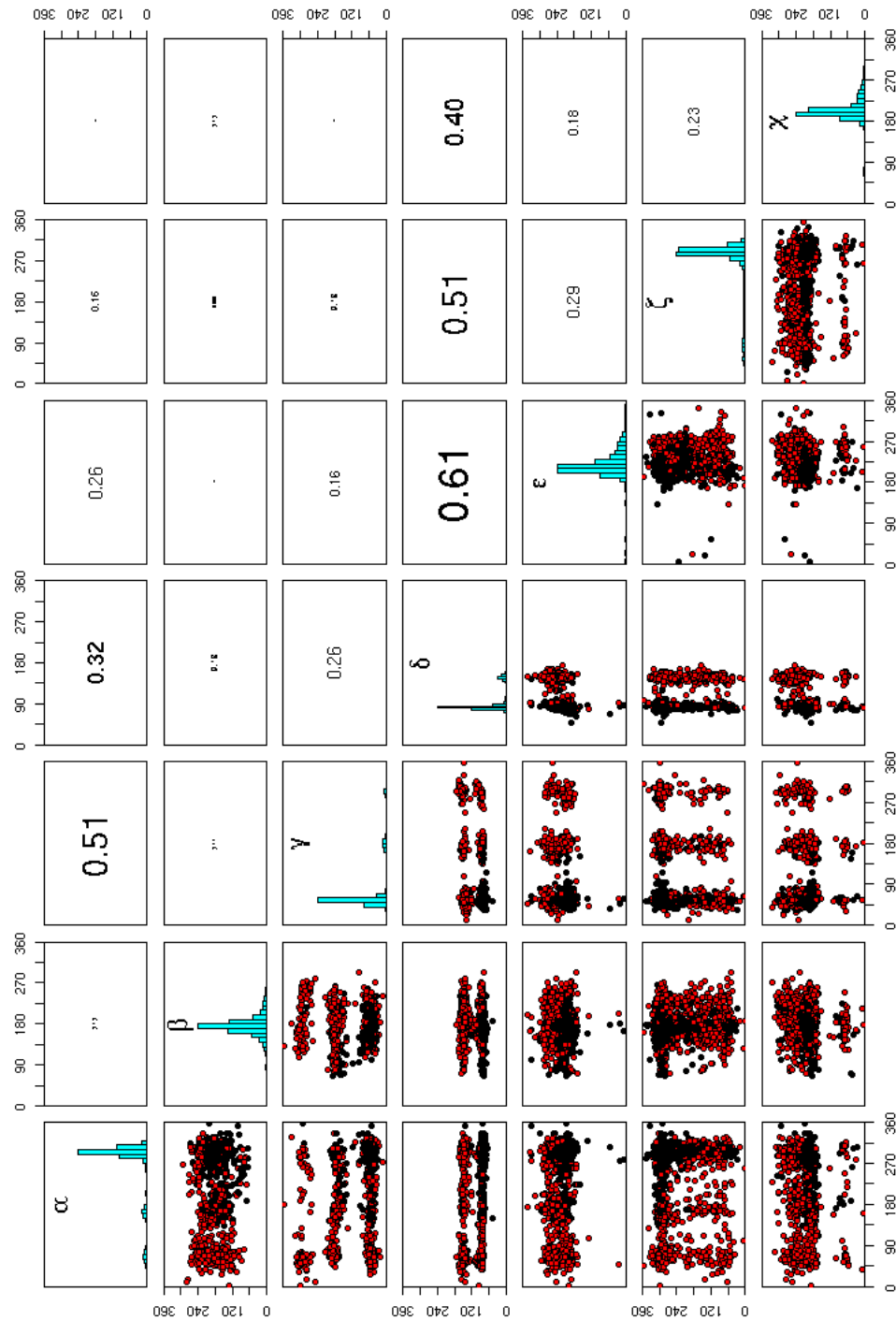


Figure 2.12: K-means of torsion angle vectors of 2753 dinucleotide steps present in 23S rRNA using the *McQueen* algorithm. The number of partitions is 2. The upper diagonal matrix displays the values of the linear correlation coefficient  $r$ , and a histogram showing the torsion angle distribution is rendered in the diagonal.

## References

- [1] Olson, W. K. and Flory, P. J. (1972) Spatial Configurations of Polynucleotide Chains. I. Steric Interactions in Polyribonucleotides: A Virtual Bond Model. *Biopolymers*, **11**, 1–23.
- [2] Saenger, W. (1984) Principles of Nucleic Acid Structure, Springer-Verlag, London.
- [3] Gautheret, D., Major, F., and Cedergren, R. (1993) Modeling the Three-dimensional Structure of RNA Using Discrete Nucleotide Conformational Sets. *Journal of Molecular Biology*, **229**(4), 1049–1064.
- [4] Wimberly, B. T., Brodersen, D. E., Clemons, W. M., Morgan-Warren, R. J., Carter, A. P., Vonnrhein, C., Hartschk, T., and Ramakrishnan, V. (2000) Structure of the 30S Ribosomal Subunit. *Nature*, **407**, 327–339.
- [5] Schlutzen, F., Tocilj, A., Zarivach, R., Harms, J., Gluehmann, M., Janell, D., Bashan, A., Bartels, H., Agmon, I., Franceschi, F., and Yonath, A. (2000) Structure of Functionally Activated Small Ribosomal Subunit at 3.3 Å Resolution. *Cell*, **102**, 615–623.
- [6] Ban, N., Nissen, P., Hansen, J., Moore, P. B., and Steitz, T. A. (August, 2000) The Complete Atomic Structure of the Large Ribosomal Subunit at 2.4 Å Resolution. *Science*, **289**(5481), 905–920.
- [7] Noller, H. F. (2005) RNA Structure: Reading the Ribosome. *Science*, **309**, 1508–1514.
- [8] Weinberg, Z., Perreault, J., Meyer, M. M., and Breaker, R. R. (2009) Exceptional Structured Non-coding RNAs Revealed by Bacterial Metagenome Analysis. *Nature*, **462**, 656–659.
- [9] Reijmers, T. H., Wehrens, R., and Buydens, L. M. C. (2001) The Influence of Different Structure Representations on the Clustering of an RNA Nucleotides Data Set. *Journal of Chemical Information and Computer Science*, **41**, 1388–1394.
- [10] Sykes, M. T. and Levitt, M. (2005) Describing RNA Structure by Libraries of Clustered Nucleotide Doublets. *Journal of Molecular Biology*, **351**, 26–38.
- [11] Murray, L. J. W., III, W. B. A., Richardson, D. C., and Richardson, J. S. (2003) RNA Backbone is Rotameric. *Proceedings of the National Academy of Sciences of the United States of America*, **100**, 13904–13909.
- [12] Hershkovitz, E., Tannenbaum, E., Howerton, S. B., Sheth, A., Tannenbaum, A., and Williams, L. D. (2003) Automated Identification of RNA Conformational Motifs: Theory and Application to the HM LSU 23S rRNA. *Nucleic Acids Research*, **31**, 6249–6257.
- [13] Schneider, B., Moravek, Z., and Berman, H. (2004) RNA Conformational Classes. *Nucleic Acids Research*, **32**, 1666–1677.
- [14] Hershkovitz, E., Sapiro, G., Tannenbaum, A., and Williams, L. D. (2006) Statistical Analysis of RNA Backbone. *Transactions on Computational Biology and Bioinformatics*, **3**, 33–46.
- [15] Duarte, C. M. and Pyle, A. M. (1998) Stepping Through an RNA Structure: A Novel Approach to Conformational Analysis. *Journal of Molecular Biology*, **284**, 1465–1478.

- [16] Duarte, C. M., Wadley, L. M., and Pyle, A. M. (2003) RNA Structure Comparison, Motif Search and Discovery Using a Reduced Representation of RNA Conformational Space. *Nucleic Acids Research*, **31**(16), 4755–4761.
- [17] Wadley, L. M., Keating, K. S., Duarte, C. M., and Pyle, A. M. (2007) Evaluating and Learning from RNA Pseudotorsional Space: Quantitative Validation of a Reduced Representation for RNA Structure. *Journal of Molecular Biology*, **372**, 942–957.
- [18] Westhof, E. and Fritsch, V. (2000) RNA folding: beyond Watson-Crick pairs. *Structure*, **8**, R55–R65.
- [19] Leontis, N. B., Stombaugh, J., and Westhof, E. (2002) The Non-Watson-Crick Base Pairs and their Associated Isostericity Matrices. *Nucleic Acids Research*, **30**, 3497–3531.
- [20] Leontis, N. B., Lescoute, A., and Westhof, E. (2006) The Building Blocks and Motifs of RNA Architecture. *Current Opinion in Structural Biology*, **16**, 279–287.
- [21] Restrepo, G., Mesa, H., Llanos, E. J., and Villaveces, J. L. (2004) Topological Study of the Periodic System. *Journal of Chemical Information and Computer Science*, **44**, 68–75.
- [22] Restrepo, G., Llanos, E. J., and Meza, H. (2006) Topological Space of the Chemical Elements and its Properties. *Journal of Mathematical Chemistry*, **39**, 401–416.

## **Chapter 3**

### **RNA Base-Pairing**

The RNA base-pairs are reviewed again.

#### **3.1 Canonical and Noncanonical Base-pairs, Methods Paper**

#### **3.2 Clustering of Yurong's Classification**

## **Chapter 4**

### **RNA Base Pair Steps**

- 4.1 Analysis (Albany Poster) and Django Webserver**
- 4.2 Persistence Length vs. Hagerman**
- 4.3 AMBER: Persistence Length of Base-Pair Step Patterns**

## Chapter 5

### RNA Motifs

#### 5.1 GNRA tetraloop

In order to compare our work to that of others on RNA structural motif localization and discovery, we ask the following questions:

1. Can the geometric rigid-block description of base-pairing and base-stacking solve the problem of defining RNA structural motifs?
2. Can we use quantities derived from the 3DNA software package to make an automatic search for a known motif, for example, the GNRA tetraloop motif, and perhaps find unknown motifs?

In the ROC meeting of May, 2009 a reduced dataset of RNA structures found at:

[http://docs.google.com/Doc?id=dhrmkfmrn\\_13ftpbjcgq](http://docs.google.com/Doc?id=dhrmkfmrn_13ftpbjcgq)

was made available to participants with the purpose of allowing them to search for RNA motifs, which would later be compared between groups. We have modestly, and as of yet unsuccessfully, started to aim at solving question number two. Initially we are trying to identify all instances of the well known GNRA tetraloop motif in the 23S subunit of ribosomal RNA of *Thermus Thermophilus*, PDB-ID:1ffk using results from 3DNA and 3DNA-Parser, and using an automated process which could be later reproduced for any desired dataset. Our hope is that these baby steps will allow us to tackle the whole ROC dataset.

##### 5.1.1 3DNA-Parser

We started by using Dr. Yurong Xin's 3DNA-Parser hoping that the description of the enclosing base pair in the loop, that is, the sheared G-A, would have a characteristic signature. We found that such is not the case. We know from Major et al. [4] that there should be at least 21 GNRA tetraloops in the 23S subunit of rRNA. We used the G2696 N2697 R2698 A2699 tetraloop as a seed (as can be seen in Figure 1.1) and found out that according to Dr. Xin's helical classification the enclosing G is classified as  $S_{hq}$  and A is classified as  $H_e$ . We then searched all such instances for G-A base-pairs and we found seven hits, but none were in fact GNRA tetraloops.

##### 5.1.2 Overlap Scores

We clustered the overlap values imposing a cutoff of values of [1-8]. There are many values which are exactly zero (33%), so, without the cutoff the zero values "overshadow" the data. For this case we obtained a "good" dendrogram as seen in Figure 1.2.

The next step in this analysis will be to find the structures which correspond to this clusters and superimpose and align them using Kabsh's algorithm to be able to determine their RMSD's.

Many people start their RNA Motif identification and classification algorithms by splitting RNA structures into what is helical and what is not, and then finding interactions between these two groups. We believe that we could do a similar exercise with 3DNA by using the scalar product of helical axis vectors



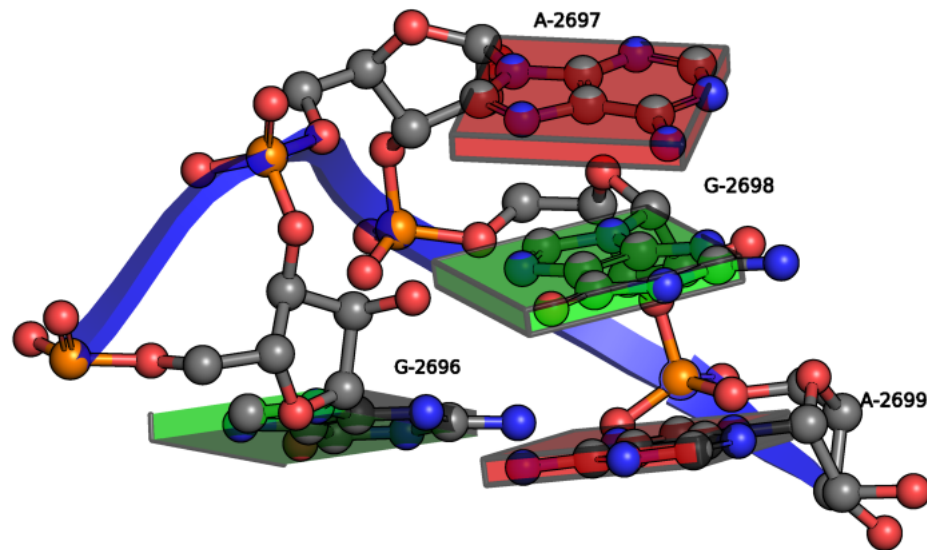


Figure 5.1: GNRA Tetraloop from *Thermus Thermophilus* 23S Ribosomal RNA PDB-ID:1ffk.

and once helical and non-helical regions are found we might be able to use 3DNA Parser to look for characteristic interactions.

## 5.2 Triplets on RNA (comparison to Laing et al.)

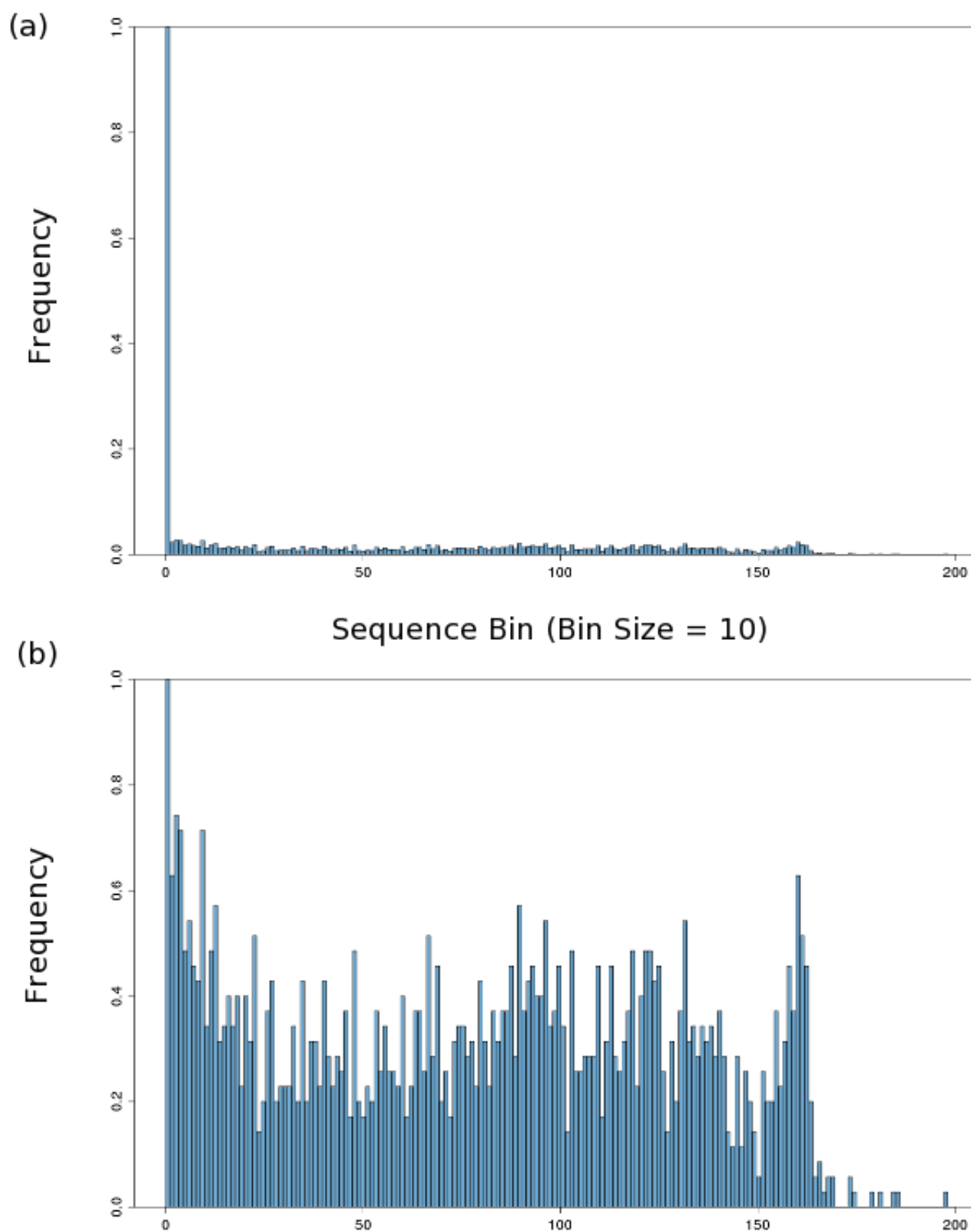


Figure 5.2: Normalized histograms showing the distribution of overlap values in the 23S subunit or *Thermus Thermophilus* rRNA, PDB-ID:1jjk. In histogram (a) all values are included, but in histogram (b) only values greater than zero are included. Notice the high preponderance of zero values, exactly 897 out of a total of 2705.

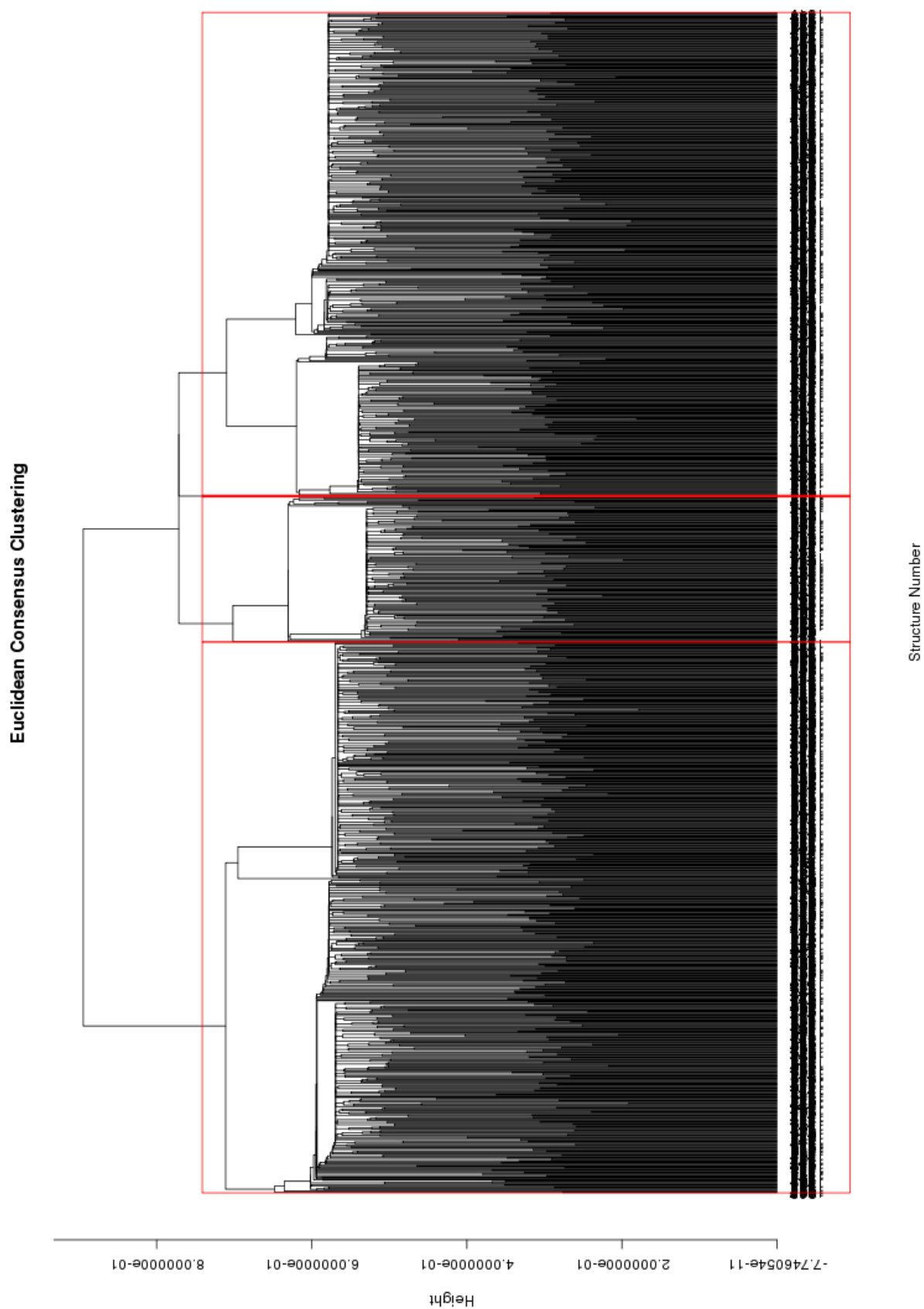


Figure 5.3: Dendrogram for consensus clustering of overlap scores in the ribosome. Zero values filtered out and remaining data normalized.

## References

- [1] Holbrook, S. R. (2005) RNA Structure: The Long and the Short of it. *Current Opinion in Structural Biology*, **15**, 302–308.
- [2] Leontis, N. B. and Westhof, E. (2003) Analysis of RNA Motifs. *Current Opinion in Structural Biology*, **13**, 300–308.
- [3] Moore, P. B. (1999) Structural Motifs in RNA. *Annual Review of Biochemistry*, **68**, 287–300.
- [4] Lemieux, S. and Major, F. (2006) Automated Extraction and Classification of RNA Tertiary Structure Cyclic Motifs. *Nucleic Acids Research*, **34**, 2340–2346.

## **Chapter 6**

### **RNA Helical Regions and Graph Theory**

Chapter on RNA Helical Region Recognition and description using graph theoretical descriptors.

## Appendix A

### Clustering Analysis (CA)

#### A.1 Hierarchical methods

The hierarchical clustering methods used were:

1. *Single linkage clustering*, where the minimum distance between elements of each cluster is taken as clustering criteria.

$$D(X, Y) = \min\{d(x_i, y_j) : x_i \in X, y_j \in Y\} \quad (\text{A.1})$$

where  $X$  and  $Y$  are vectors, and  $d(x_i, y_j)$  is the distance between cluster elements.

2. *Complete linkage clustering*, where the maximum distance between cluster elements is the clustering criteria.

$$D(X, Y) = \max\{d(x_i, y_j) : x_i \in X, y_j \in Y\} \quad (\text{A.2})$$

3. *Average linkage clustering*, the mean distance between elements of each cluster is taken as clustering criteria.

$$D(X, Y) = \frac{1}{N_x * N_y} \sum_{i=1}^{N_x} \sum_{j=1}^{N_y} d(x_i, y_j) \quad (\text{A.3})$$

where  $N_x$  and  $N_y$  are the number of elements in respective clusters.

4. *Centroid linkage clustering*, uses the distance between cluster centroids, as clustering criteria.

$$D(X, Y) = d(\bar{x}, \bar{y}) \quad (\text{A.4})$$

$$\bar{x} = \frac{1}{N_x} \sum_{i=1}^{N_x} x_i \quad (\text{A.5})$$

$$\bar{y} = \frac{1}{N_y} \sum_{i=1}^{N_y} y_i \quad (\text{A.6})$$

$$(\text{A.7})$$

Structure	Property I	Property II
1	1.00	5.00
2	-2.00	6.00
3	2.00	-2.00
4	-2.00	-3.00
5	3.00	-4.00

Table A.1: Example of structures, considered as bidimensional vectors, to be clustered using the average linkage method and the Manhattan distance.

5. *Ward's Method*, uses the error sum of squares (ESS).

$$D(X, Y) = ESS(XY) - [ESS(X) + ESS(Y)] \quad (\text{A.8})$$

$$ESS(X) = \sum_{i=1}^{N_x} \left| x_i - \frac{1}{N_x} \sum_{j=1}^{N_x} x_j \right|^2 \quad (\text{A.9})$$

As an example lets think of a case where we have five structures. Each one of them is described by a bidimensional vector as illustrated in Table A.1.

The first step is to chose a distance definition. We chose Manhattan and the distance values between structures can be displayed in a lower triangular matrix as seen in equation A.10

$$d(X, Y) = \begin{vmatrix} & 1 & 2 & 3 & 4 \\ 1 & & & & \\ 2 & 4 & & & \\ 3 & 8 & 12 & & \\ 4 & 11 & 9 & 5 & \\ 5 & 11 & 15 & 3 & 6 \end{vmatrix} \quad (\text{A.10})$$

Let's calculate explicitly the Manhattan distance between structures 2 and 3,

$$d(2, 3) = |-2.00 - 6.00| + |2.00 - -2.00| = 12 \quad (\text{A.11})$$

Now that we have calculated the distances we need a clustering method, in this case, we will use the average linkage clustering method. The first step is to group whatever structures are closer, that is, structures 3 and 5 ( $d(3, 5) = 3$ ). Now we find the mean distance between the elements of this cluster and the remaining unclustered structures, that is, structures 1, 2 and 4, we obtain the following mean distances

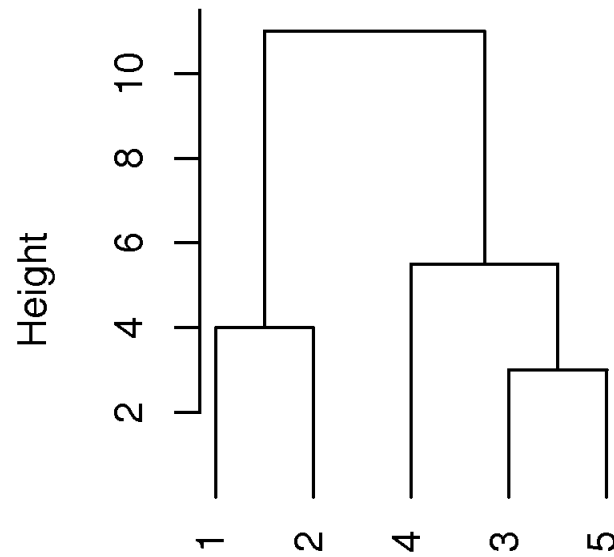
$$D(\{3, 5\}, 1) = \frac{1}{2 * 1} * (8 + 11) = 4.5 \quad (\text{A.12})$$

$$D(\{3, 5\}, 2) = \frac{1}{2 * 1} * (12 + 15) = 13.5 \quad (\text{A.13})$$

$$D(\{3, 5\}, 4) = \frac{1}{2 * 1} * (5 + 6) = 5.5 \quad (\text{A.14})$$

Since the distances between  $\{3, 5\}$  and all remaining unclustered vectors is higher than the distance between vectors 1 and 2 ( $d(1, 2) = 4$ ) then  $\{1, 2\}$  are grouped. The following value, in hierarchical increasing order is 4.5 between  $\{3, 5\}$  and 1 (see equation A.12), but since 1 and 2 are already grouped we can't group  $\{3, 5\}$  with 1. The next value, following the lower to higher hierarchy, is 5 ( $d(3, 4) = 5$ ),

## Average linkage example tree



## Manhattan distance

Figure A.1: Clustering tree for 5 bidimensional vectors using the Manhattan distance definition and the average linkage clustering method.

but we have already grouped 3 with 5, so we have to keep advancing in the hierarchy. The next value is 5.5, which corresponds to grouping  $\{3, 5\}$  with 4, so we cluster them. The only remaining possibility for grouping is, group  $\{1, 2\}$  and  $\{4, 3, 5\}$ , so we do it as illustrated in Figure A.1.



## Appendix B

### Dimension Reduction

#### B.1 Principal Component Analysis

Given a set of data it's important to check before analysis if the dimensions of such set can be reduced. A very common method to check this is called Principal Component Analysis (PCA).

The PCA method can be defined in a strictly mathematical way as the method which; finds the Principal Components (PCs) of a dataset by doing "an orthogonal linear transformation of a set of variables optimizing certain algebraic criterion" [1].

$$\mathbf{y} = \mathbf{T}\mathbf{x} \quad (\text{B.1})$$

Where  $\mathbf{T}$  is an orthogonal linear transformation matrix of dimension  $k$  by  $n$ ,  $\mathbf{x}$  is the "original" data matrix of dimension  $n$  by  $m$ , and by definition of the matrix product  $\mathbf{y}$  is a matrix of dimension  $k$  by  $m$ .

From the linear tranformation expression it's clear that if  $k = m$  then the transformation matrix is just a rotation matrix, and in the case where  $k < m$  then the transformation matrix is also reducing the dimension of the "original" data.

One common algorithm to find such transformation ( $\mathbf{T}$ ) is the following:

1. Substract the mean from the data matrix  $\mathbf{x}$ .

$$\mathbf{\Omega} = \mathbf{x} - \bar{\mathbf{x}}_m \quad (\text{B.2})$$

2. Find the covariance matrix for  $\mathbf{\Omega}$ .

$$\mathbf{\Sigma} = \frac{\mathbf{\Omega}^{-1}\mathbf{\Omega}}{(1 - n)} \quad (\text{B.3})$$

3. Diagonalize the covariance matrix  $\mathbf{\Sigma}$ .

$$\mathbf{T}^T \mathbf{\Sigma} \mathbf{T} = \mathbf{\Lambda} \quad (\text{B.4})$$

4. Organize  $\mathbf{T}$  from the highest to lowest eigenvalues in  $\mathbf{\Lambda}$ .

Obtaining the eigenvalues and eigenvectors of  $\mathbf{\Sigma}$  means that we have found our transformation matrix  $\mathbf{T}$ , which can be used to either rotate the original data space to an orthonormal one, or to reduce the dimensionality of the data space by choosing  $k < m$ , depending on the weight of the eigenvectors in  $\mathbf{\Lambda}$ . The  $k$  rows of  $\mathbf{y}$  are named Principal Components.

There is a large amount of bibliography which refers to the statistics of Principal Components Analysis, and to sofistications which are not included in this brief appendix. An interested reader can find great help in the following web addresses:

<http://www-stat.wharton.upenn.edu/~buja/script-Buja-CU-2009-06-pca.R>

## References

- [1] Jolliffe, I. T. (2002) Principal Component Analyses, Springer, .

## **Appendix C**

### **Figure Supplements**

#### **C.1 Chapter2**

These are additional figures for chapter 2.

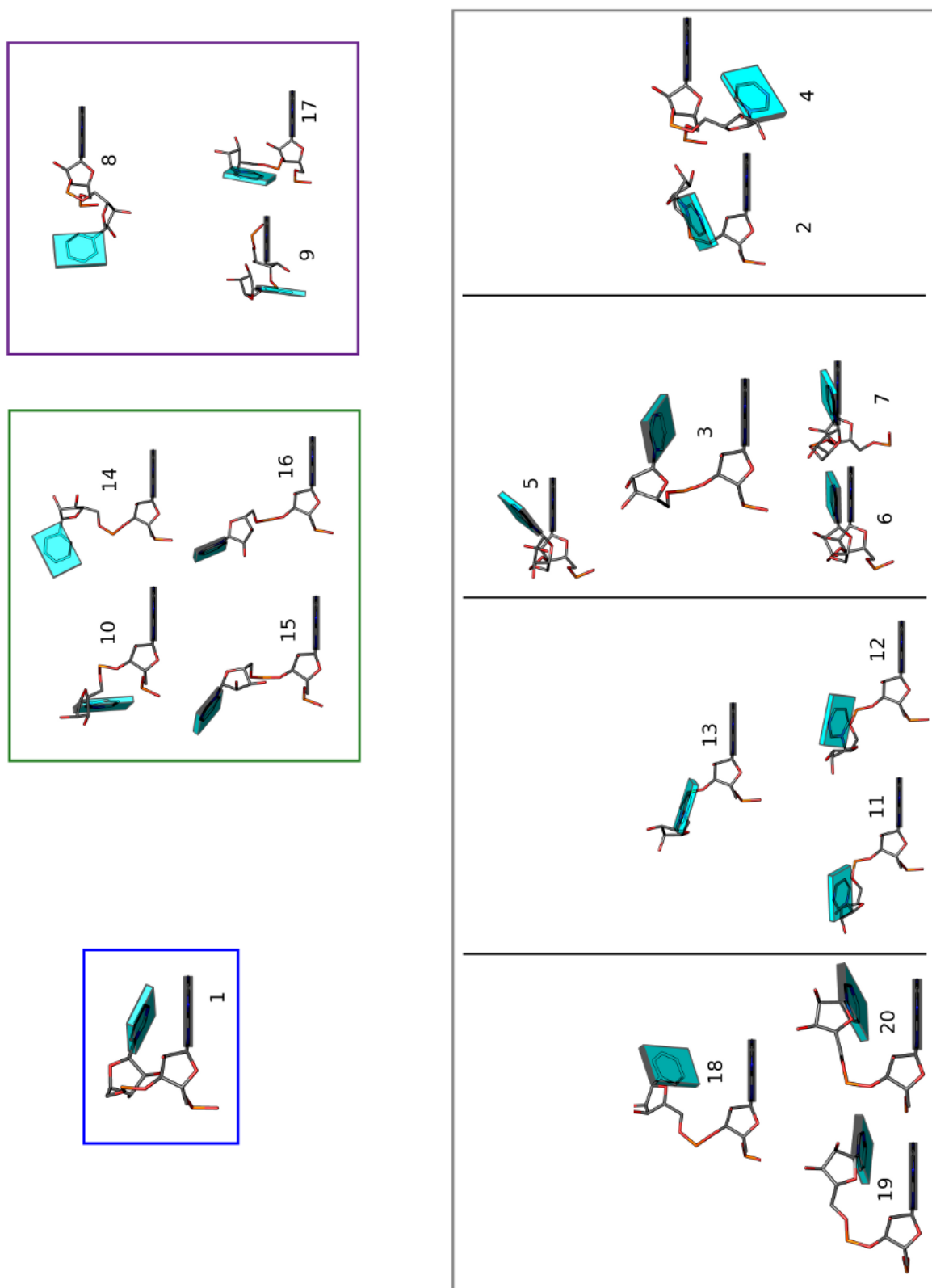


Figure C.1: Non A-RNA Type base steps centered on the standard reference frame of Adenine. Top view with the Minor Groove side of Adenine pointing down the page and the Major Groove pointing up.

## Curriculum Vitae

### Mauricio Esguerra

#### La Mala Educacion

- 1991** High School Diploma from Gimnasio Moderno, Bogota, Colombia.
- 2000** B. Sc. in Chemistry from Universidad Nacional de Colombia
- 2010** Ph. D. in Chemistry and Chemical Biology, Rutgers University

#### Professional Experience

- 2003-2009** Teaching assistant, Department of Chemistry and Chemical Biology, Rutgers University

#### Publications

- 2009** W. K. Olson, M. Esguerra, Y. Xin, X-J. Lu, Methods