# RNA Comparative Structure Modeling...
## three steps ahead



**Marc A. Marti-Renom**
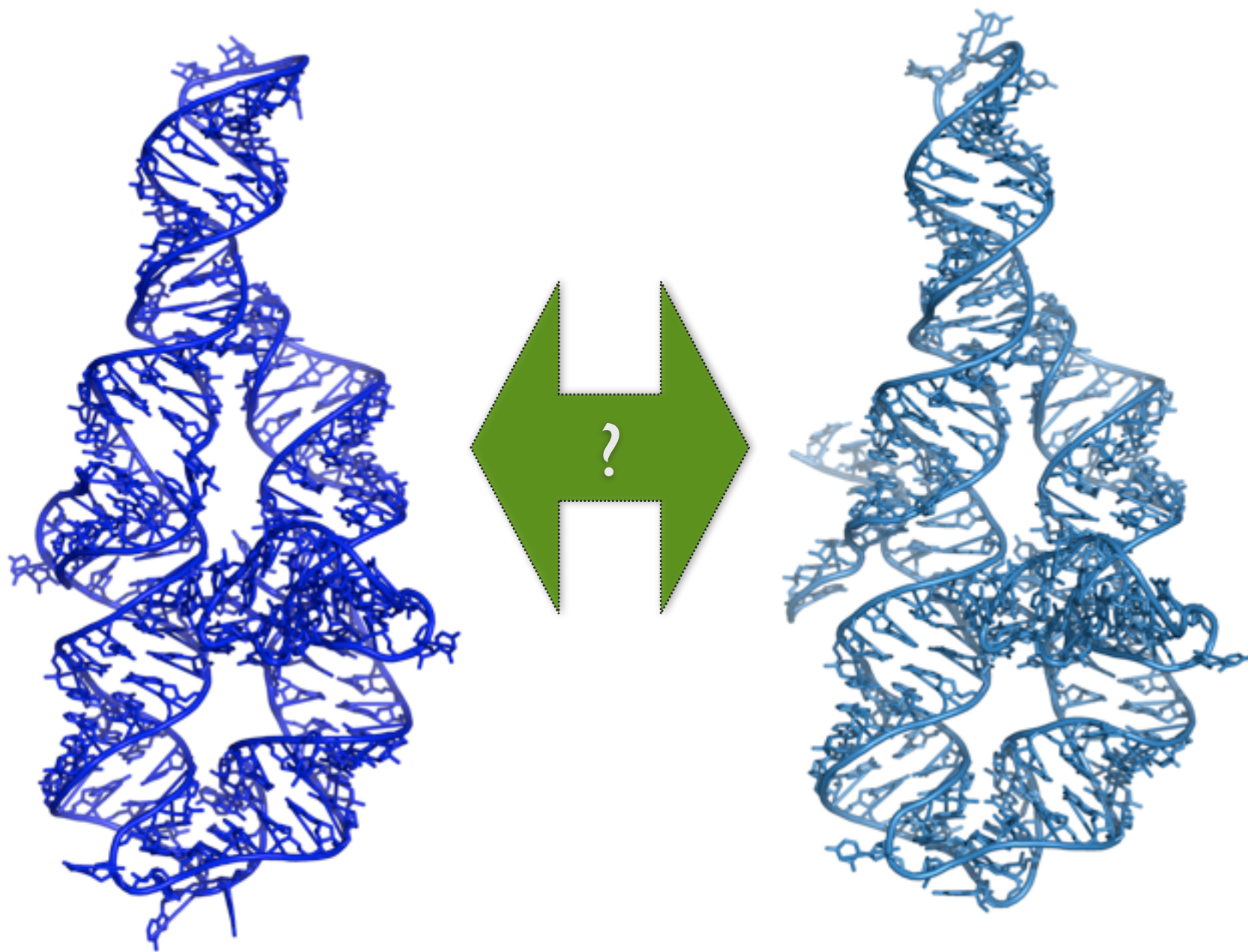http://sgu.bioinfo.cipf.es

Structural Genomics Unit
Bioinformatics & Genomics Department
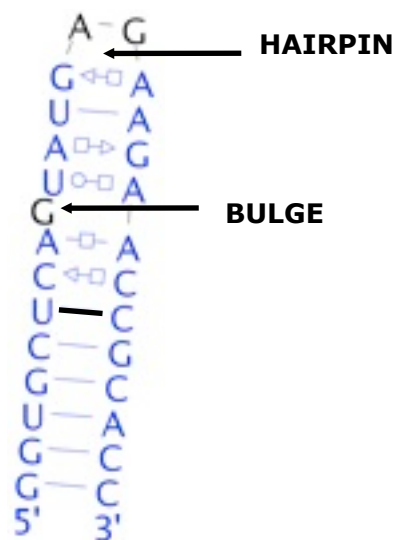Prince Felipe Resarch Center (CIPF), Valencia, Spain

# First step

## Can we reliably compare RNA structures?

# RNA structure

### Primary Structure

>Mutant Rat 28S rRNA sarcin/ricin domain
GGUGCUCAGUAUGAGAAGAACCGCACC
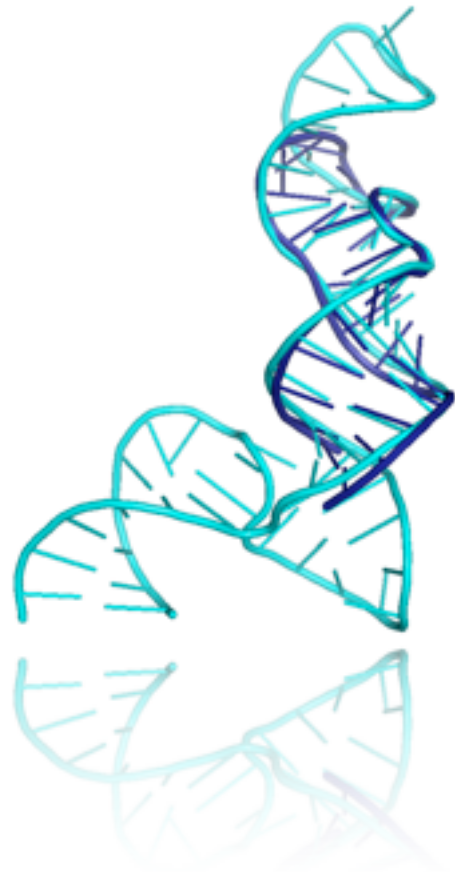
### Secondary Structure

>Mutant Rat 28S rRNA sarcin/ricin domain
GGUGCUCAGUAUGAGAAGAACCGCACC
( ( ( ( ( ( ( ( ( . ( ( ( ( . . ) ) ) ) ) ) ) ) ) ) ) )

### Tertiary Structure

Secondary Structure interactions and other interactions like pseudoknots, hairpin-hairpin interactions etc.

3

# Structural alignment

Structural alignment attempts to establish equivalences between two or more polymer structures based on their shape and three-dimensional conformation.
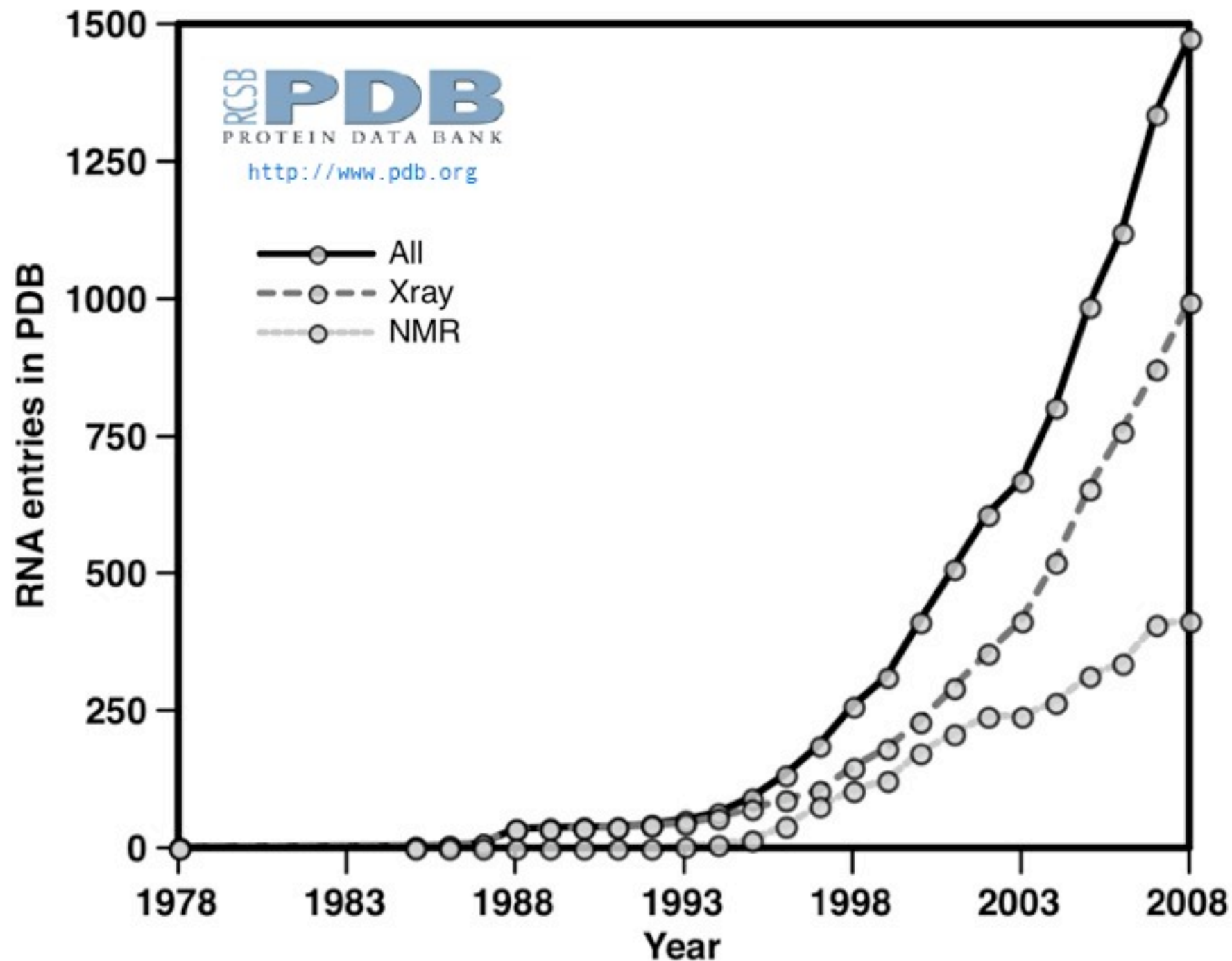
In contrast to simple structural superposition, where at least some equivalent residues of the two structures are known, structural alignment does not require prior knowledge of the equivalent positions.

Structural alignment has been used as a valuable tool for the comparison of proteins, including the inference of evolutionary relationships between proteins of remote sequence similarity.

4

# RNA Structure

Currently **more than 1500 RNA structures** are deposited in the PDB (Mar 09)

# RNA structure datasets

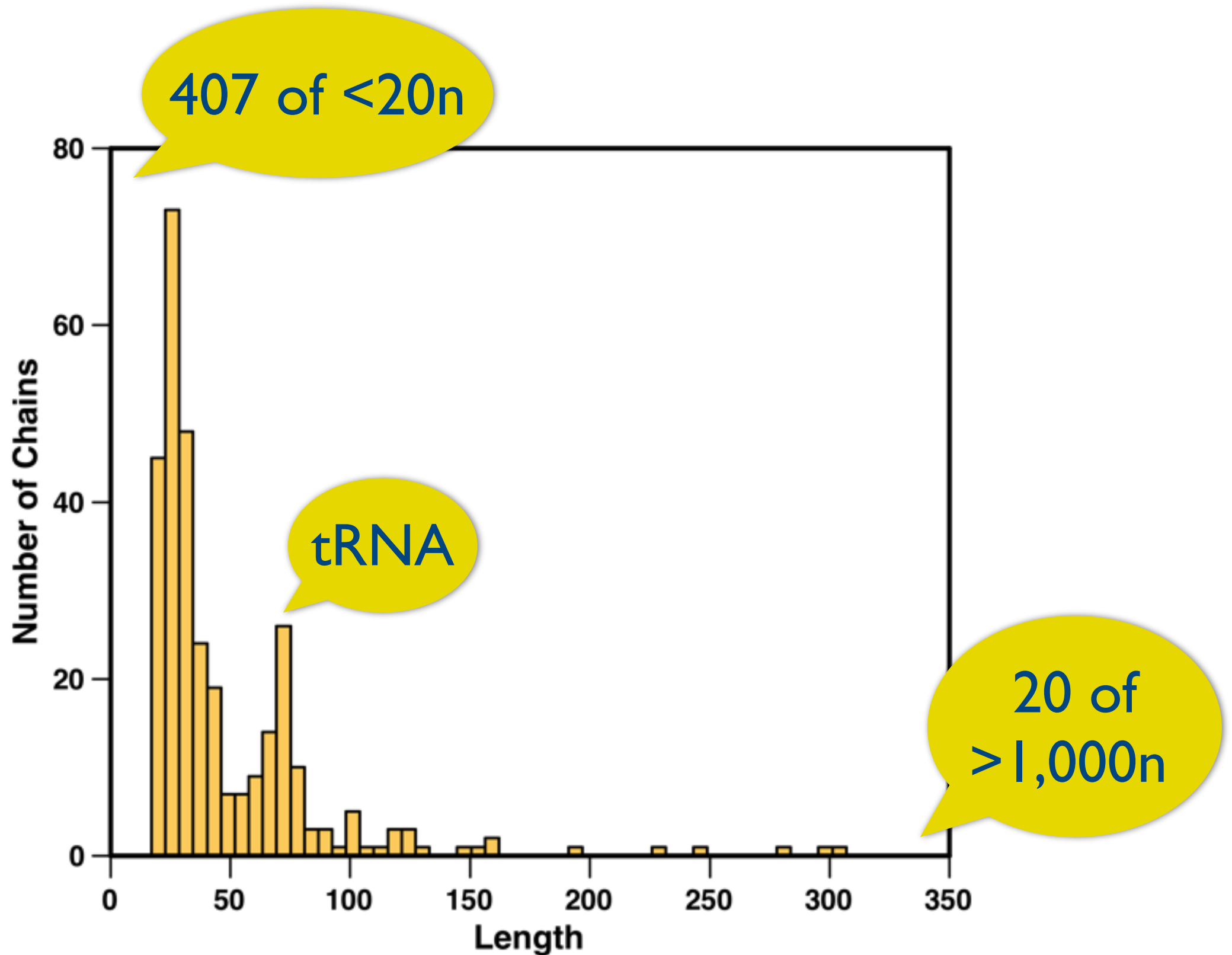| RNA STRUCTURE* | 1,101 |
|---|---|
| RNA CHAINS | 2,179 |
| Non-Redundant RNA CHAINS** | 744 |
| RNA CHAINS (20≤ Length ≤310) | 313 |
| HIGH RESOLUTION RNA SET*** | 54 |

NR95

HR

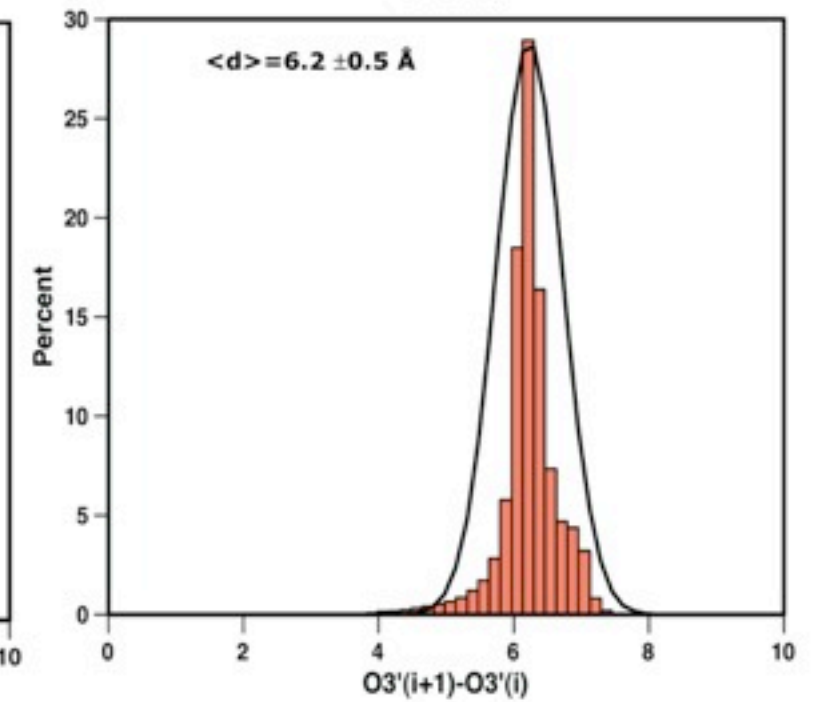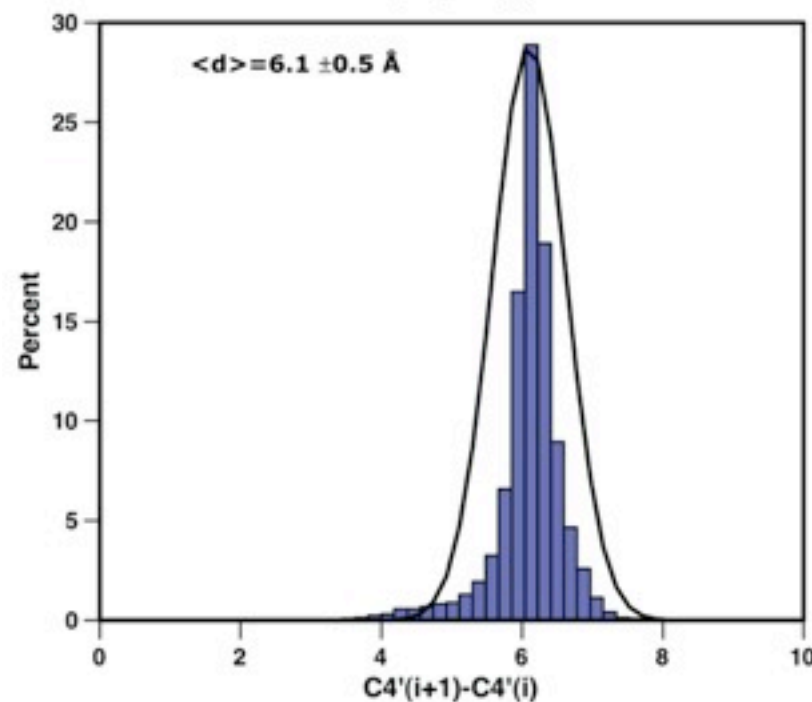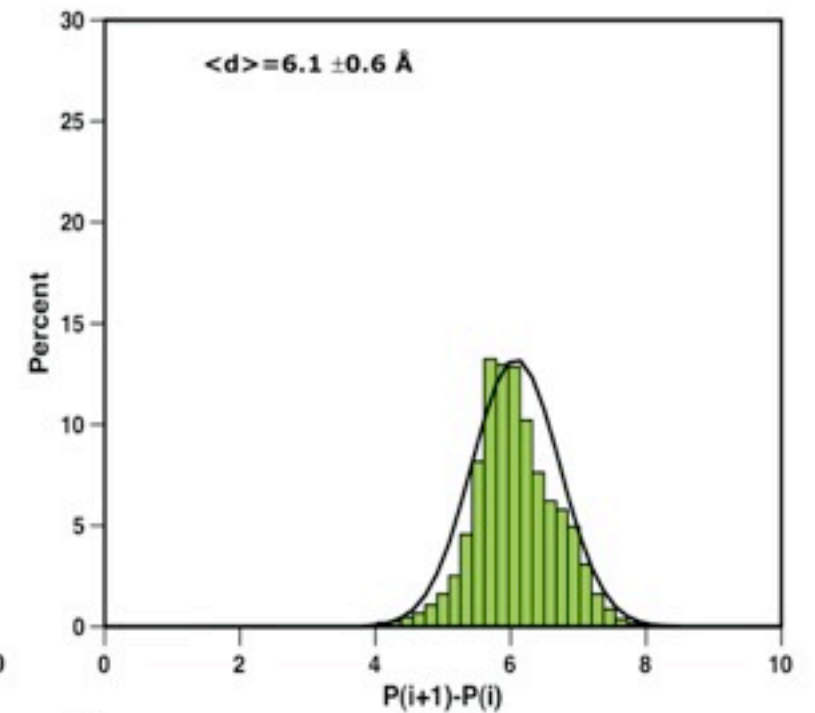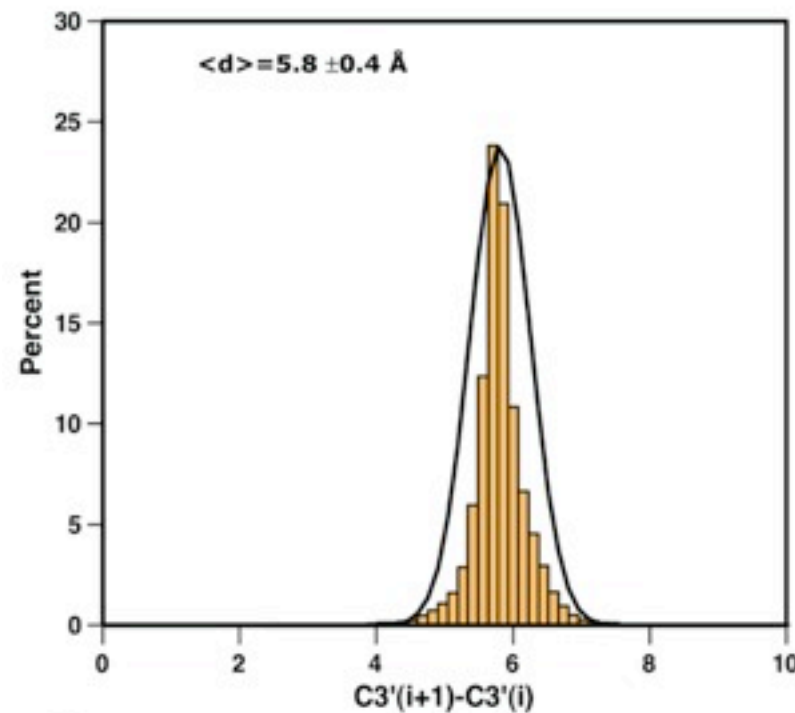\* *from PDB November 06.*

\*\* *non-redundant 95% sequence identity*
\*\*\* *Resolution below 4.0 Å and with no missing backbone atoms.*
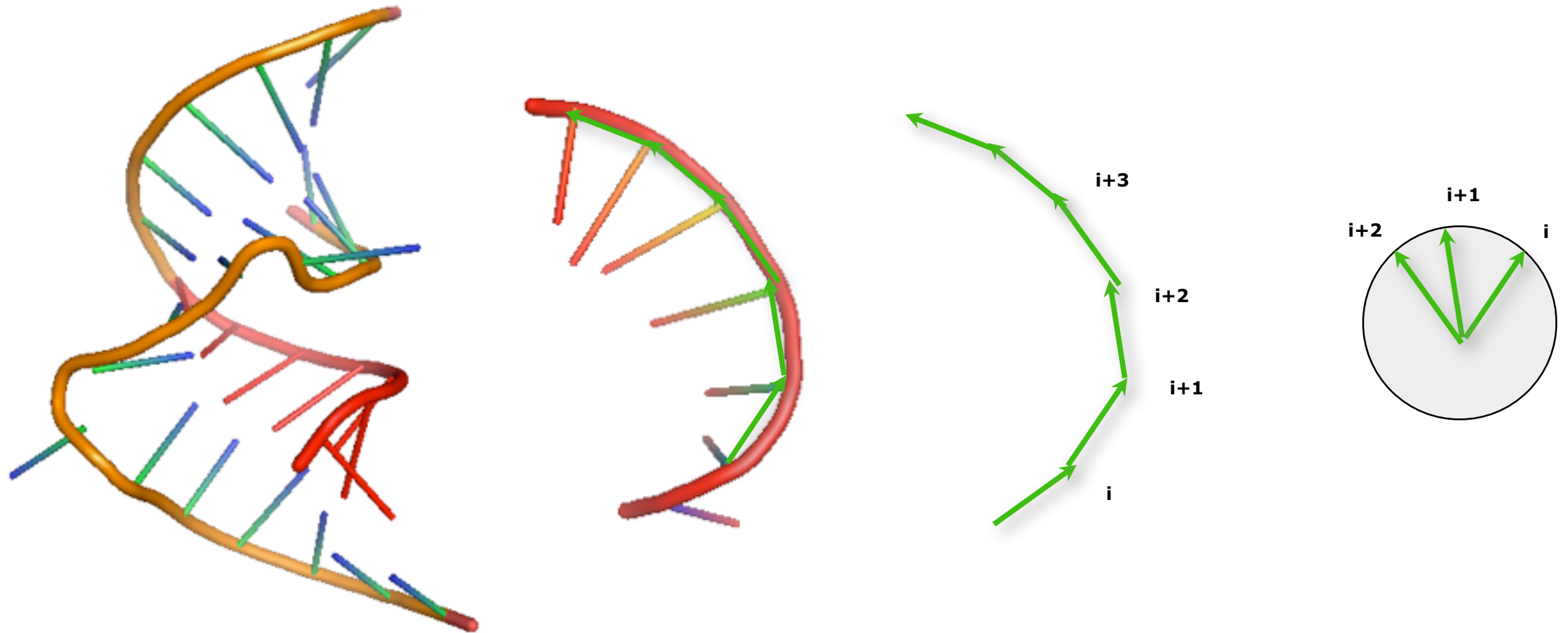
6

# Dataset distribution

# Atom selection

The best backbone atom that represents the RNA structure has been selected by evaluating the distribution of the distances between consecutive atoms in structures from the NR95 set.

# Unit Vector I

## Representation



A Unit Vector is the normalized vector between two successive C3' atoms.

For each position *i* consider the *k* consecutive vectors, which will be mapped into a unit sphere representing the local structure of k residues.

# Unit Vector II

## Scoring



$$URMS^R = \sqrt{2.0 - \frac{2.84}{\sqrt{k}}}$$

$$S_{ij} = \frac{(URMS^R - URMS^{ij})}{URMS^R} \Delta(URMS^R, URMS^{ij})$$

$$\Delta(URMS^R, URMS^{ij}) = 10 \Rightarrow URMS^R > URMS^{ij}$$

$$\Delta(URMS^R, URMS^{ij}) = 0 \Rightarrow URMS^R \leq URMS^{ij}$$

For each position i, the k consecutive unit vectors are grouped and aligned to the j set of unit vectors. Each pair of aligned unit vectors will be evaluated by calculating Unit Root Mean Square distance ($URMS^{ij}$).

The obtained URMS values are compared the minimum expected URMS distance between two random set of k unit vectors ($URMS^R$).

The alignment score is than calculated normalizing $URMS^{ij}$ to the $URMS^R$ value.

10

# Alignment

Sq/St 1  1 ———●——— N i
Sq/St 2  1 ———●——— M j



$$D_{i,j}=\min \begin{cases} D_{i,j-1}+\text{Score}_{(Ä,rj)} \\ D_{i-1,j-1}+\text{Score}_{(ri,rj)} \\ D_{i-1,j}+\text{Score}_{(ri,Ä)} \end{cases}$$

Best alignment score

Backtracking to get the best alignment

A Dynamic Programming procedure is then applied to search for the optimal structural alignment using a global alignment with zero end gap penalties.

The maximum subset of local structures that have their corresponding C3' within 3.5 Å in the space are evaluated. The number of close atoms is used to evaluate the percentage of structural identity (PSI) using a variant of the MaxSub algorithm.

Needleman and Wunsch J. Mol.Biol 1970
Siew et al. Bioinformatics 2000

Thursday, July 23, 2009

# Random RNA

In order to build a background distribution that reproduce the scores given by the structural alignments of unrelated RNA sequences, we generated a set 300 random RNA sequences and structures with sequence length uniformly distributed between 20 and 320 nucleotides.

The RNA backbone can be described given the 6 torsion angle ($\alpha,\beta,\gamma,\delta,\epsilon,\zeta$) for each nucleotide.

The RNA backbone is rotameric and only 42 conformation have been described from a set o high resolution structures .

According to this observation we generated the 300 structures, randomly selecting the backbone angles among the 42 possible conformations.

Murray et al PNAS 2003

# Background distribution

Considering a dataset of 300 random RNA structures, we have produced ~45,000 pairwise alignments that resulted in a empirical distribution. From such distribution we can then evaluate μ and σ needed to calculated the p-value for P(s>=x).

Empirical



Analytic



$$P(s \geq x) = 1 - \exp(-e^{-\lambda(s-\mu)})$$

Karlin and Altschul, 1990 *PNAS* **87**, pp2264

Thursday, July 23, 2009

# Mean and sigma

The score distribution depends on the length of the molecule.

We divided the resulting structural alignments (~45,000) in 30 bins according to the shorter sequence length of the two random structures (N).

For each bin the μ and σ values are evaluated fitting the data to an EVD.

The relations between N and μ, σ values are extrapolate fitting them to a power low function (r≈0.99).



Legend:
$\mu = 763 * \mathbf{N}^{-0.896}$
$\sigma = 180 * \mathbf{N}^{-1.010}$

x-axis: N (Length of the shorter RNA structure)

# Statistical significance

all-against-all comparison of structures in the NR95 set

# SARA .vs. ARTS



**SARA**

Percentage of structure identity (PSI)   92.6%
Percentage of sequence identity   48.0%
Percentage of SSE identity 100.0%
RMSD 1.78 Å

```
>1q96 Chain:A
-------------------ggugcucaguaugag--------aagaaccgcacc-------
>1un6 Chain:E
gccggccacaccuacggggccugguuaguaccugggaaaccugggaauaccaggugccggc
```

**ARTS**

Percentage of structure identity (PSI)   76.9%
Percentage of sequence identity   20.0%
Percentage of SSE identity 79.2%
RMSD 1.66Å

```
>1q96 Chain:A
-------------------gugcucaguaugaga-----aga-accgcacc--------
>1un6 Chain:E
ccggccacaccuacggggccugguuaguaccugggaaaccugggaauaccaggugccggc
```

**PSI:**  % of structure identity

**PSS:** % of secondary structure identity

**Cut-off distance:** 4.0 Å

16

# SARA Alignments



A) Staphylococcus phage group I ribozyme (1y0qA)
   Human group I Intron fragment (1u6bB)

| Aligned nucleotides: | 120 |
|---|---|
| RMSD: | 1.8 Å |
| Sequence Identity: | 34.0 % |
| Secondary Structure Identity: | 52.1 % |
| Structure Identity: | 60.9 % |
| Sequence -ln(p-value): | 18.2 |
| Secondary structure -ln(p-value): | 10.3 |
| Structure -ln(p-value): | 15.6 |
| **Mean -ln(p-value):** | **14.7** |

B) Pyrococcus horikoshii tRNA(Leu) (1wz2C)
   Acuifex aeolicus tRNA(Met) (2ct8C)

| Aligned nucleotides: | 65 |
|---|---|
| RMSD: | 1.9 Å |
| Sequence Identity: | 56.8 % |
| Secondary Structure Identity: | 88.5 % |
| Structure Identity: | 87.8 % |
| Sequence -ln(p-value): | 10.2 |
| Secondary structure -ln(p-value): | 5.2 |
| Structure -ln(p-value): | 7.2 |
| **Mean -ln(p-value):** | **7.5** |

C) Synthetic P4-P6 RNA ribozyme (1l8vA)
   Mus musculus P4-P6 RNA ribozyme (2r8sR)

| Aligned nucleotides: | 134 |
|---|---|
| RMSD: | 1.8 Å |
| Sequence Identity: | 80.9 % |
| Secondary Structure Identity: | 81.0 % |
| Structure Identity: | 85.4 % |
| Sequence -ln(p-value): | 37.0 |
| Secondary structure -ln(p-value): | 17.1 |
| Structure -ln(p-value): | 19.4 |
| **Mean -ln(p-value):** | **24.5** |

D) Haloarcula marismortui 23S RNA (3cce0)
   Thermus thermophilus 23S RNA (3d5bA)

| Aligned nucleotides: | 2,347 |
|---|---|
| RMSD: | 1.7 Å |
| Sequence Identity: | 52.7 % |
| Secondary Structure Identity: | 75.7 % |
| Structure Identity: | 85.2 % |
| Sequence -ln(p-value): | 37.0 |
| Secondary structure -ln(p-value): | 37.0 |
| Structure -ln(p-value): | 37.0 |
| **Mean -ln(p-value):** | **37.0** |

17

# Second step...

## Can we reliably predict RNA function from structure?

# RNA function annotation



*Capriotti and Marti-Renom Bioinformatics 2008*
*Tamura et al. NAR 2004*

19

# Results

| Datasets | Number of chains | Number of alignments | Number of different SCOR functions |
|----------|------------------|----------------------|-------------------------------------|
| RNA09 | 451 | 101 475 | |
| BgALI | 451 | 50 995 | |
| FSCOR | 419 | | 168 |
| R-FSCOR | 192 | | 168 |
| T-FSCOR | 227 | | 88 |

**leave one out on FSCOR**

**performances on T-FSCOR**



20

# Third step...

To what extend can we do comparative RNA structure prediction?

21

# Stx/Seq relationship



$$PSI = 97.6\,(1 - e^{-0.051 * PID})$$

$$R = 0.71$$

$$p = 8.7 \cdot 10^{-103}$$

# SSE/Stx/Seq relationship

# Twilight Zone

# SARA server

Capriotti et al. Bioinformatics (2008) vol. 24 (16) pp. i112-i118

25

# SARA server

E. Capriotti, M. A. Marti-Renom (2008), *Bioinformatics* **24**:i112          E. Capriotti, M. A. Marti-Renom. (2009). *NAR* **37**:W260

# Acknowledgments

http://sgu.bioinfo.cipf.es