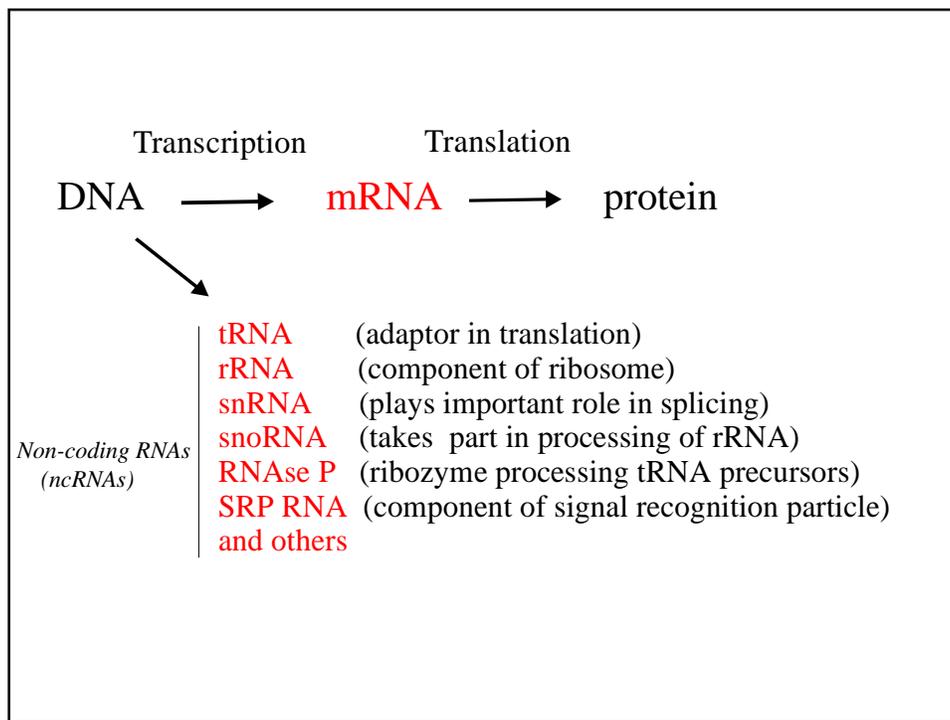


# RNA bioinformatics

Two major problems are addressed:

- \* How do we find RNA genes computationally?
- \* How do we predict the secondary structure of RNAs?

Tore Samuelsson Oct 2005



## *The RNA world*

Early in the evolution of life RNA molecules played an important role, and they were to a large extent independent of protein and DNA.

RNA molecules have two important properties to make this possible:

- \* they are able to carry genetic information
- \* they are able to catalyse chemical reactions

Many RNAs today are “molecular fossils” from the RNA world.

For instance, a number of important cellular processes such as

*splicing &  
protein synthesis*

are catalyzed by RNA molecules (ribozymes)

RNA molecules, unlike DNA molecules,  
are single-stranded

**DNA**

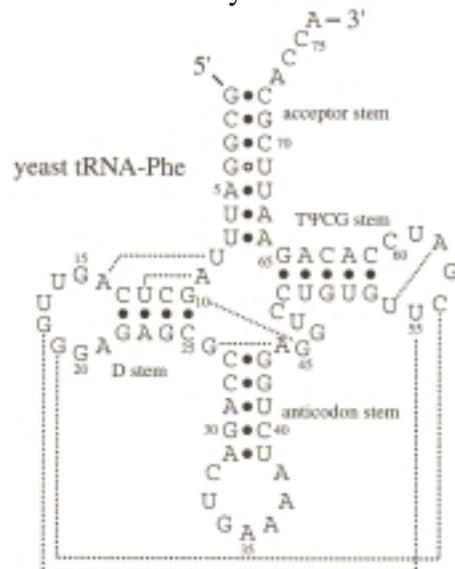
5' GCCAAGGTTTCGAAA 3'  
3' CGGTTCCAAGCTTT 5'

**RNA**

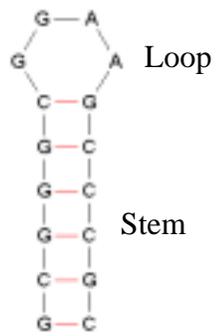
5' GCCAAGGUUCGAAA 3'

... but through internal base-pairing  
helical structures occur like those in the  
double-stranded DNA .

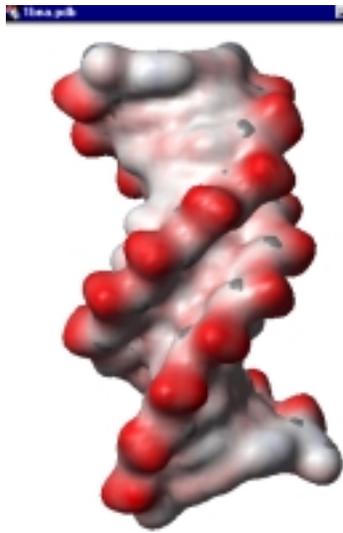
Each family of ncRNA typically adopts a  
characteristic secondary structure



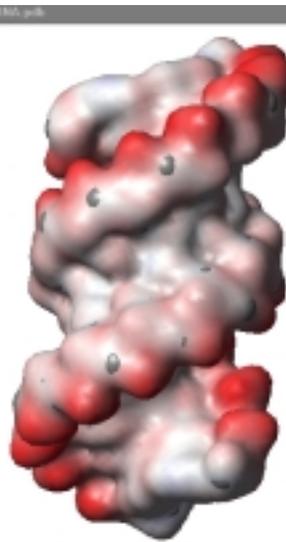
## Hairpin

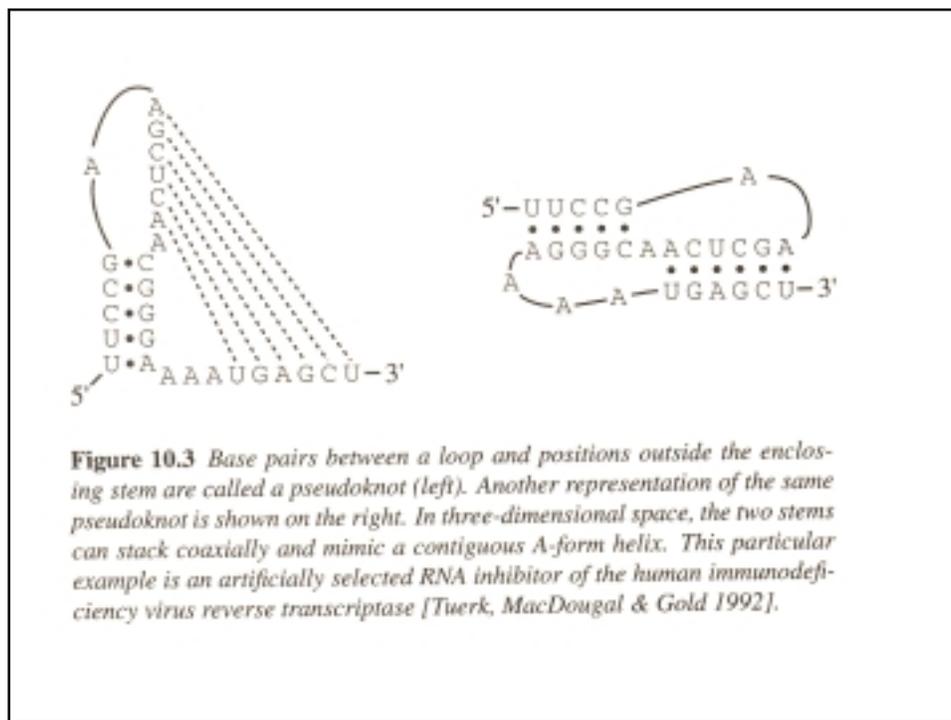
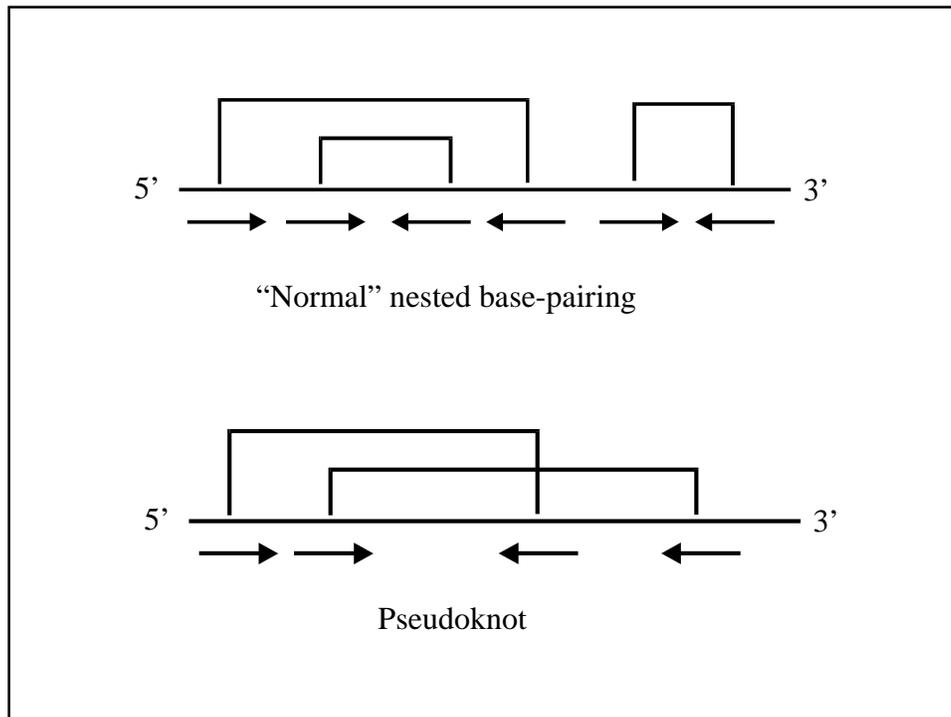


## DNA

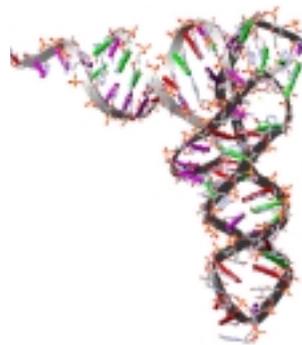
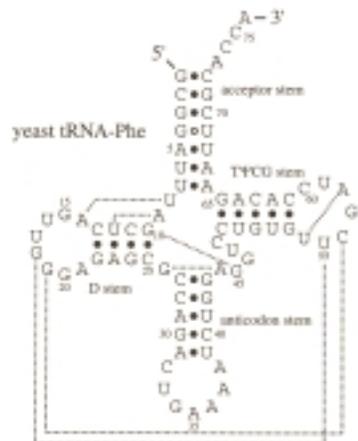
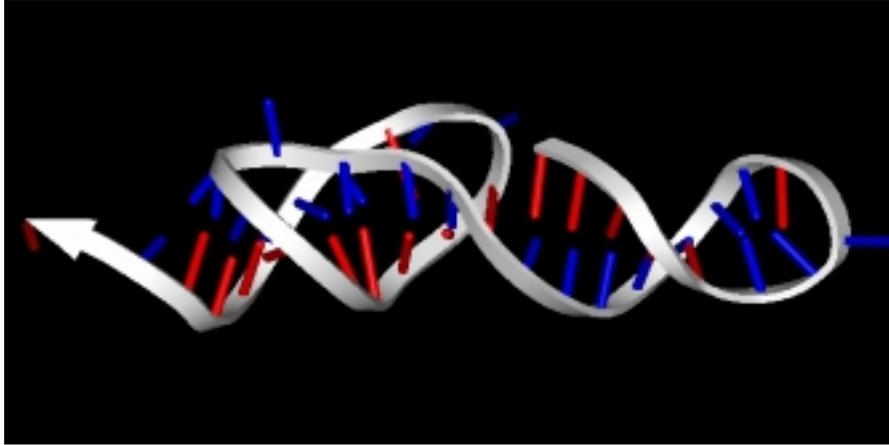


## RNA

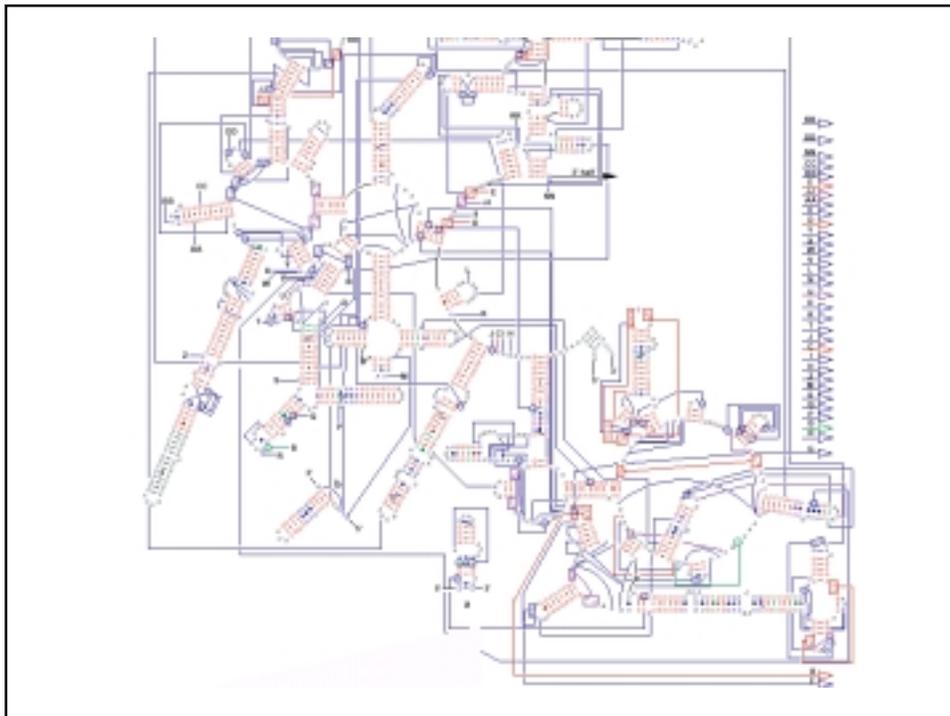
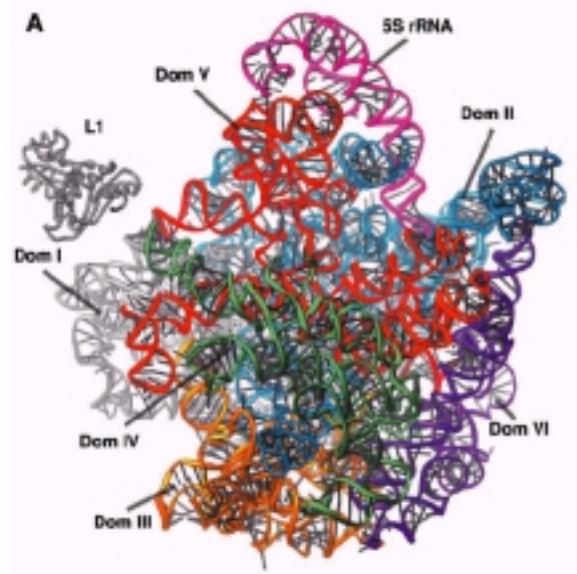




### Pseudoknot 3D structure

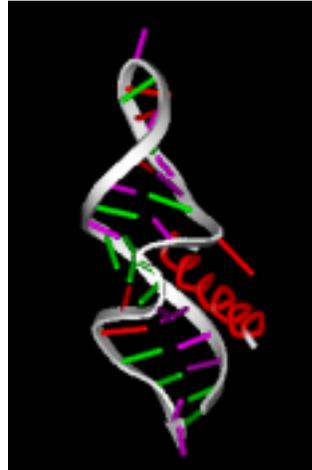
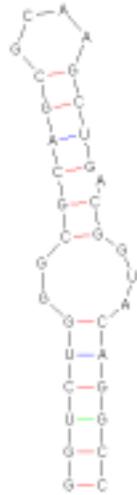


50S subunit of ribosome  
RNA 2900 nucleotides





## RNA secondary structure as a regulator of gene expression in HIV - The RRE / Rev complex



## How do we identify RNA genes ?

### Predicting protein and RNA genes

#### Protein

Identifiable primary sequence signals such as promoters, coding sequences, polyadenylation sites, exon/intron signals

Methods based on homology to previously known seqs

```
A  MVAKQRIRMANEKHKNITQRGNVAKTSRNASPEEKASVGPWLLALFIFVVC
   :.  ::::  .:::  :::::  :::::  .  :  :  .  :::::  :::::
B  MAPKQRMTLANKQFSKNVNNRGNVAKSLKPA-EDKYPAAPWLIGLFVVFVC
```

#### RNA

? No identifiable primary sequence signal that is shared by all RNA genes

? Poor sequence homology

The primary sequence is poorly conserved  
in many RNAs

=>

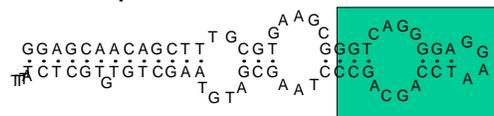
Many RNAs or functional elements in RNA  
cannot be identified by sequence alignment methods

-->

*Methods for RNA identification rely on analysis of  
both primary sequence and secondary structure.*

Example of RNA homologs with limited  
primary sequence homology -  
- bacterial SRP RNAs

**Streptococcus mutans**



**Rickettsia conorii**



Ricket	GCTAGTAGTGG.GCATTGTAC...CT...GTTTAGTCGGTCAGGT	38
Strept	GGAGCA...ACAGCTTTGCGTGAAGCGGTCAAGG	32
Ricket	CTGAAAGGAAGCAGCC...AGAGTGGGATTGATG.GGTC.ATTA	78
Strept	GAGGAATCCAGCAGCCCTAAGCGATGTAAGCTGTGTCTCTATTT	77
Ricket	CTAGCATTA	87
Strept		



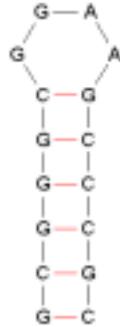
**Figure 10.4** The consensus binding site for R17 phage coat protein. N, Y and R are standard 'degenerate' symbols for multiple possible nucleotides. N indicates {A, C, G, U}, Y indicates {C, U} and R indicates {A, G}. N' indicates a complementary base pairing to N.

What methods may be used to identify RNA genes that encode RNAs that belong to known families?

- Pattern matching
- Covariance models
- Mfold (prediction of secondary structure)

## Specifying patterns in PatScan

```
p1=5...7
GGAA
~p1
```



```
r1={au,ua,gc,cg,gu,ug}
p1=6...7
GGG [1,0,0]
p2=8...9
4...5
r1~p2[1,0,1]
3...4
~p1
```

```
GGTTTGC GGA CGTATGGA
AGTG TTCACTGCG AAAA
GCAAACC
```

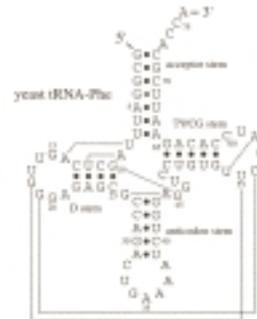
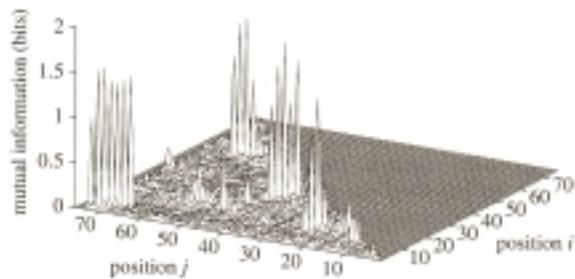




Secondary structure can be inferred by  
**comparative sequence analysis**



**Figure 10.5** Comparative sequence analysis recognises that the two boxed positions in this example of a multiple alignment (left) are covarying to maintain Watson–Crick complementarity. This covariation implies a base pair, leading to a consensus secondary structure prediction (right).



Intuitively,  $M_{ij}$  tells us how much information we get about the identity of the residue in one position if we are told the identity of the residue in the other position. In the case of a base pair with no sequence constraints, we get 2 bits

$$M_{ij} = \sum_{x_i, x_j} f_{x_i, x_j} \log_2 \frac{f_{x_i, x_j}}{f_{x_i} f_{x_j}}$$

$f_{x_i}$  is the frequency of one of the four bases (A, C, G, U) observed in column  $i$ .  $f_{x_i, x_j}$  is the joint (pairwise) frequency of one of the sixteen possible base pairs observed in columns  $i$  and  $j$ .  $M_{ij}$  measures how much the joint frequency distribution deviates from the distribution that is expected if the two columns vary independently. For the four-letter RNA alphabet,  $M_{ij}$  varies between 0 and 2 bits.  $M_{ij}$  is maximal if  $i$  and  $j$  individually appear completely random ( $f_i = f_j = 0.25$ ), but  $i$  and  $j$  are perfectly correlated, for instance in a Watson–Crick base pair.

Intuitively,  $M_{ij}$  tells us how much information we get about the identity of the residue in one position if we are told the identity of the residue in the other position. In the case of a base pair with no sequence constraints, we get 2 bits

### Mutual information example

Consider the multiple alignment:

```
GGCC
GCCG
GACU
GUCA
```

The  $M_{ij}$  expression (Durbin p. 266) gives:

Example 1)

For columns 2 and 4:

For the pair GC :  $0.25 * \log_2(0.25/(0.25*0.25)) = 0.5$

For the pair CG :  $0.25 * \log_2(0.25/(0.25*0.25)) = 0.5$

For the pair AU :  $0.25 * \log_2(0.25/(0.25*0.25)) = 0.5$

For the pair UA :  $0.25 * \log_2(0.25/(0.25*0.25)) = 0.5$

The sum gets  $0.5 * 4 = 2$  bits, which means that we have the highest possible mutual information between columns 2 and 4.

Example 2)

For columns 1 and 3:

There is only pair, GC :

$1 * \log_2(1/(1*1)) = 0$ , i.e we have no mutual information at all.

RNA genes that belong to an already known RNA family can efficiently be identified using ‘covariance models’

Rfam RNA families database of alignments and CMEs

Rfam Home Page

Rfam is a large collection of multiple sequence alignments and covariance models covering many common non-coding RNA families. For each family in Rfam you can:

- View and download multiple sequence alignments
- Read family annotation
- Examine species distribution of family members
- Follow links to other databases

Rfam makes use of a large amount of available data, especially published multiple sequence alignments, and repackages these data in a single searchable and sustainable resource. We have made every effort to credit individual sources on family pages, and have compiled a list of links to these sources [here](#). If you find any of the data presented here useful, please do be sure to credit the primary source.

For more information on Rfam, and using this site, click [here](#).

Rfam UK is also the home of the miRNA Registry. To submit miRNA sequences for naming, or to browse published miRNAs, click [here](#).

Rfam Mirror Servers Worldwide

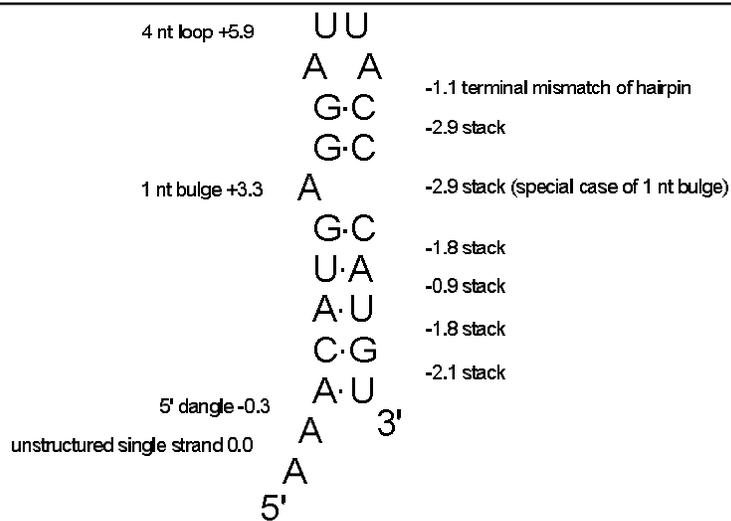
FTP access to Rfam

Name	No. seed	No. full	Average Length	Average %id	Description
<a href="#">5_SG_rRNA</a>	64	12366	151.9	76	5-8S ribosomal RNA
<a href="#">5S_rRNA</a>	606	3690	115.5	58	5S ribosomal RNA
<a href="#">6S</a>	7	18	168.6	74	6S / 5sS RNA
<a href="#">CobB</a>	6	14	351.8	75	CobB/RomB RNA family
<a href="#">Dna5</a>	2	11	86.0	90	Dna5 RNA
<a href="#">Gcd8</a>	4	7	113.0	98	Gcd8 RNA
<a href="#">Hammerhead</a>	68	305	35.9	68	Hammerhead ribozyme
<a href="#">Histone3</a>	66	808	26.0	78	Histone 3' UTR stem-loop
<a href="#">Intron_gp1</a>	30	2690	212.9	40	Group I catalytic intron
<a href="#">Intron_gp2</a>	117	2963	77.7	54	Group II catalytic intron
<a href="#">let-7</a>	12	28	79.7	69	let-7 microRNA precursor
<a href="#">McfE</a>	8	15	91.2	90	McfE RNA
<a href="#">Oxy5</a>	6	12	109.3	95	Oxy5 RNA
<a href="#">RNase_MRP</a>	25	39	265.1	52	RNase MRP
<a href="#">RNaseP_kat_a</a>	226	379	318.0	62	Bacterial RNase P class A
<a href="#">RNaseP_kat_b</a>	24	46	328.2	65	Bacterial RNase P class B
<a href="#">RNaseP_nuc</a>	47	67	290.3	41	Nuclear RNase P
<a href="#">Rps4</a>	8	10	108.2	84	Rps4 RNA
<a href="#">RRE</a>	65	486	227.8	95	HIV Rev response element
<a href="#">SECIS</a>	65	217	64.3	43	Selenocysteine insertion sequence
<a href="#">Spot_42</a>	4	12	117.9	95	Spot 42 RNA





**Mfold** predicts optimal and suboptimal secondary structures for an RNA (or DNA) molecule using the energy minimization method of Zuker



Overall  $\Delta G = -4.6$  kcal/mol

$\Delta G$  calculation for an RNA stem loop