

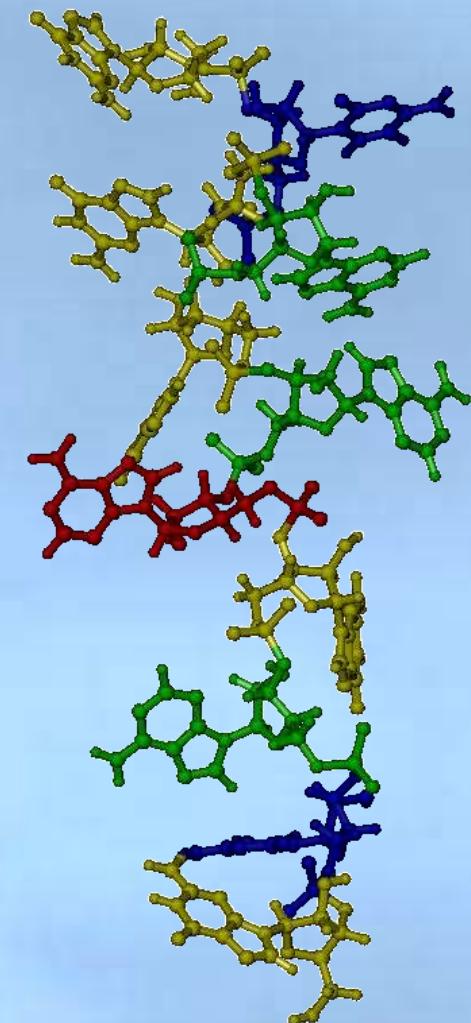
# RNA fragments database: project & implementation

---

# Agenda

---

1. Introduction
  - trends in structural bioinformatics
  - about RNA...
2. Project of RNA fragments database
  - main concept
  - database scheme & contents
3. RNA FRABASE
  - implementation
  - search examples
  - performance

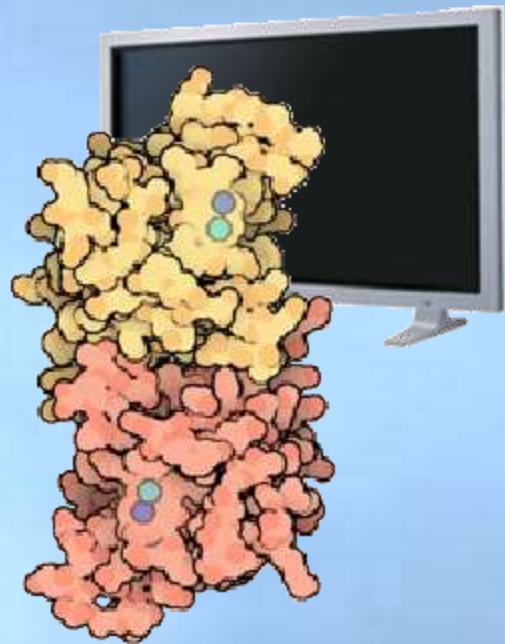


# Structural bioinformatics

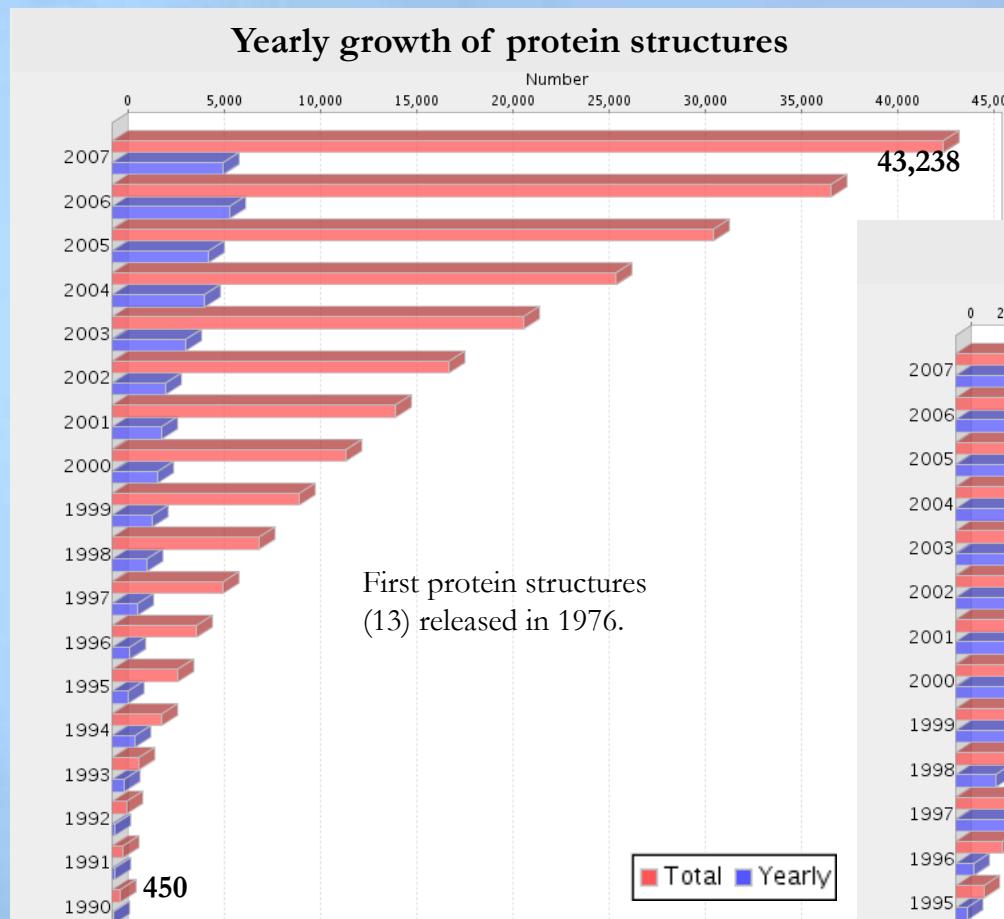
---

Basic streams in structural bioinformatics:

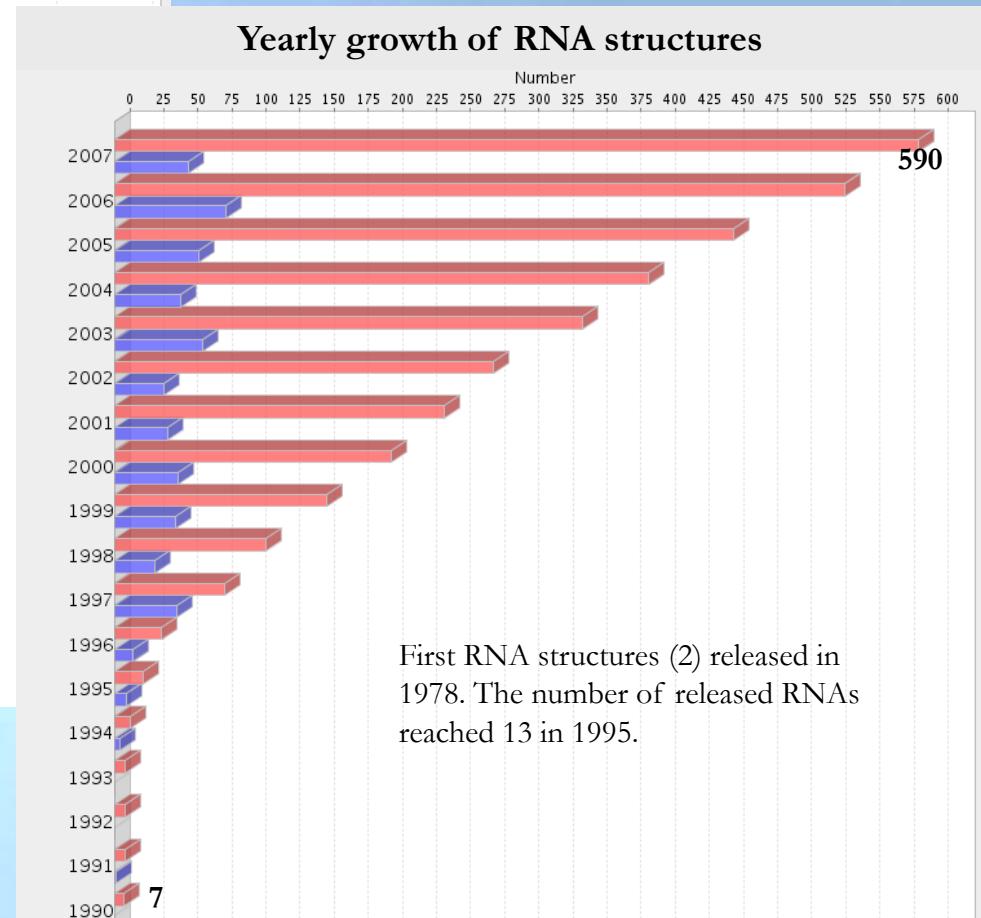
- > • Structure determination
- > • Discovering dependencies: molecular structure - functions
- > • Structure prediction
- > • Structure comparison
- ➡ • Structural databases
- Visualization of molecular models



# Trends in structural research



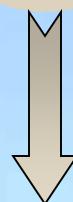
Growth of released structures per year: proteins vs RNA



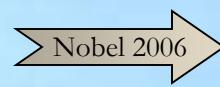
Statistics updated: Tuesday, 30.10.2007

# About RNA...

Proteins + DNA + RNA = molecular model of life



C  
G  
U  
A

- transmits genetic information from DNA into proteins,
- catalyzes chemical reactions,
- controls certain chemical processes in the cell,
- initiates the process of gene silencing  (RNAi),
- forms the genetic material of some viruses,
- is a component of the ribosomes.

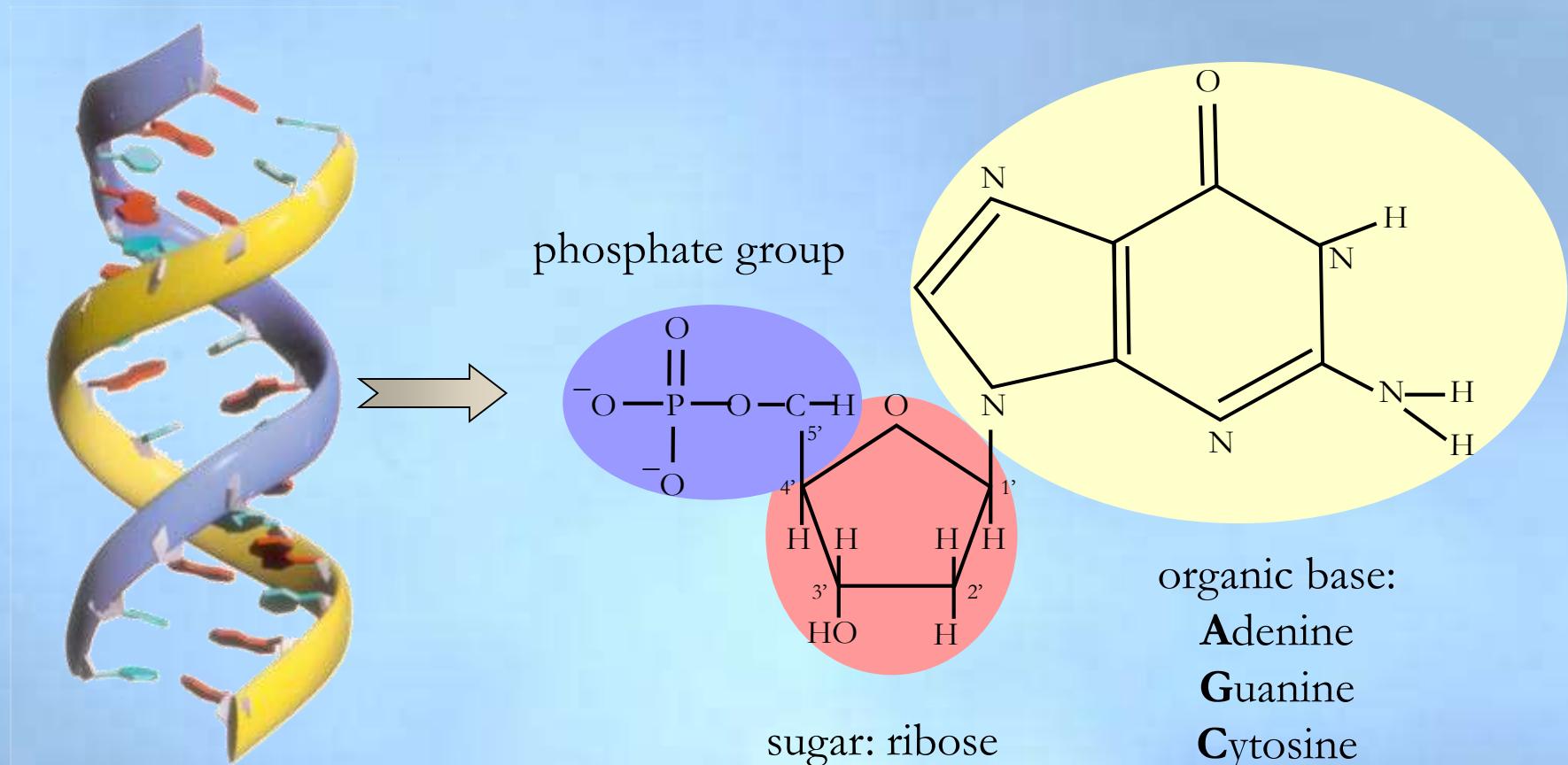


Andrew Z. Fire i Craig C. Mello - laureaci Nagrody Nobla z medycyny i fizjologii w 2006 roku.

Fot. MICHAEL PROBST AP

# RNA structure

basic component of RNA chain: nucleotide



# I- and II-dimensional structure of RNA

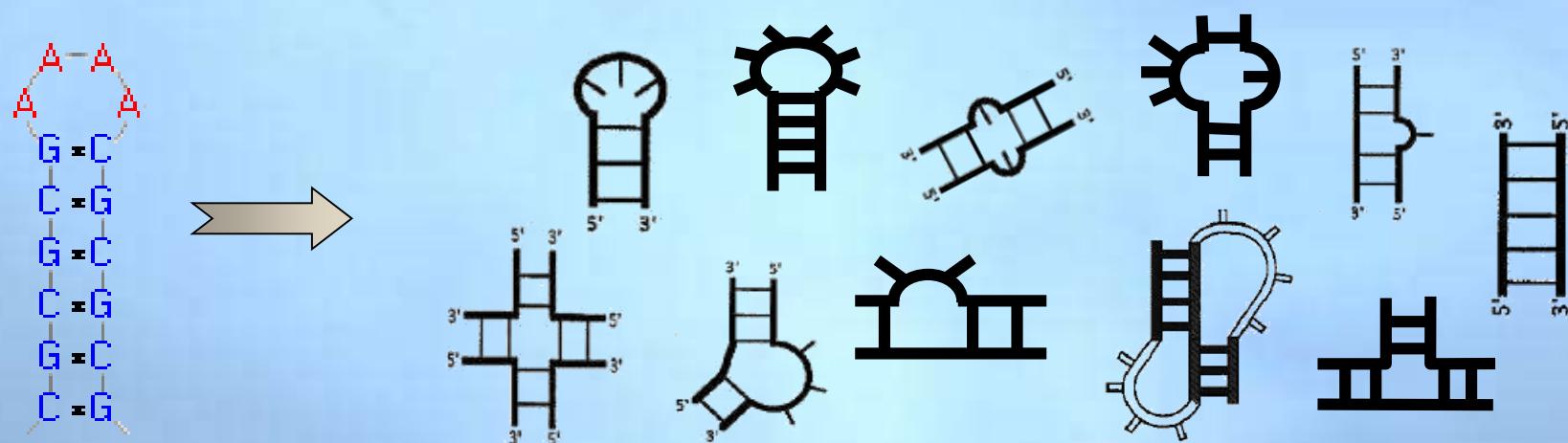
## Primary structure

describes nucleotide sequence in the chain

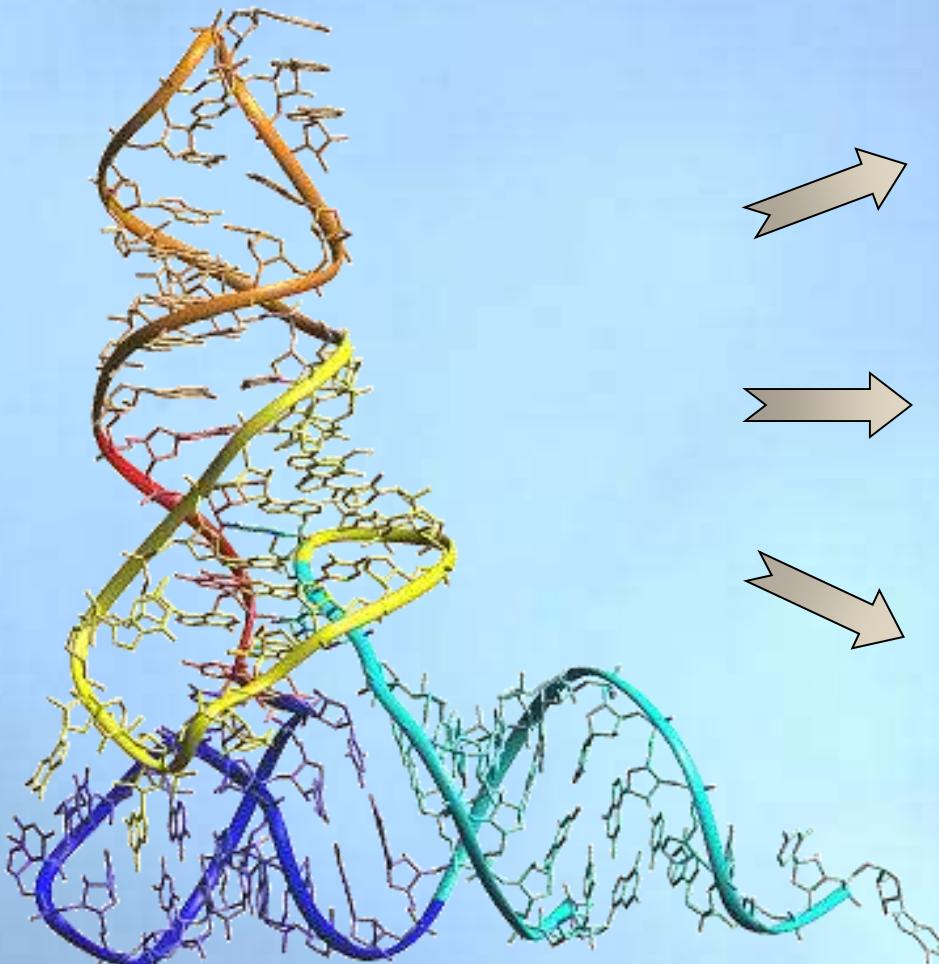
5' – C<sub>1</sub> G<sub>2</sub> C<sub>3</sub> G<sub>4</sub> A<sub>5</sub> U<sub>6</sub> C<sub>7</sub> U<sub>8</sub> G<sub>9</sub> – 3'

## Secondary structure

describes one- and two-dimensional fragments (non-/canonical base-pairs)



# III-dimensional structure of RNA

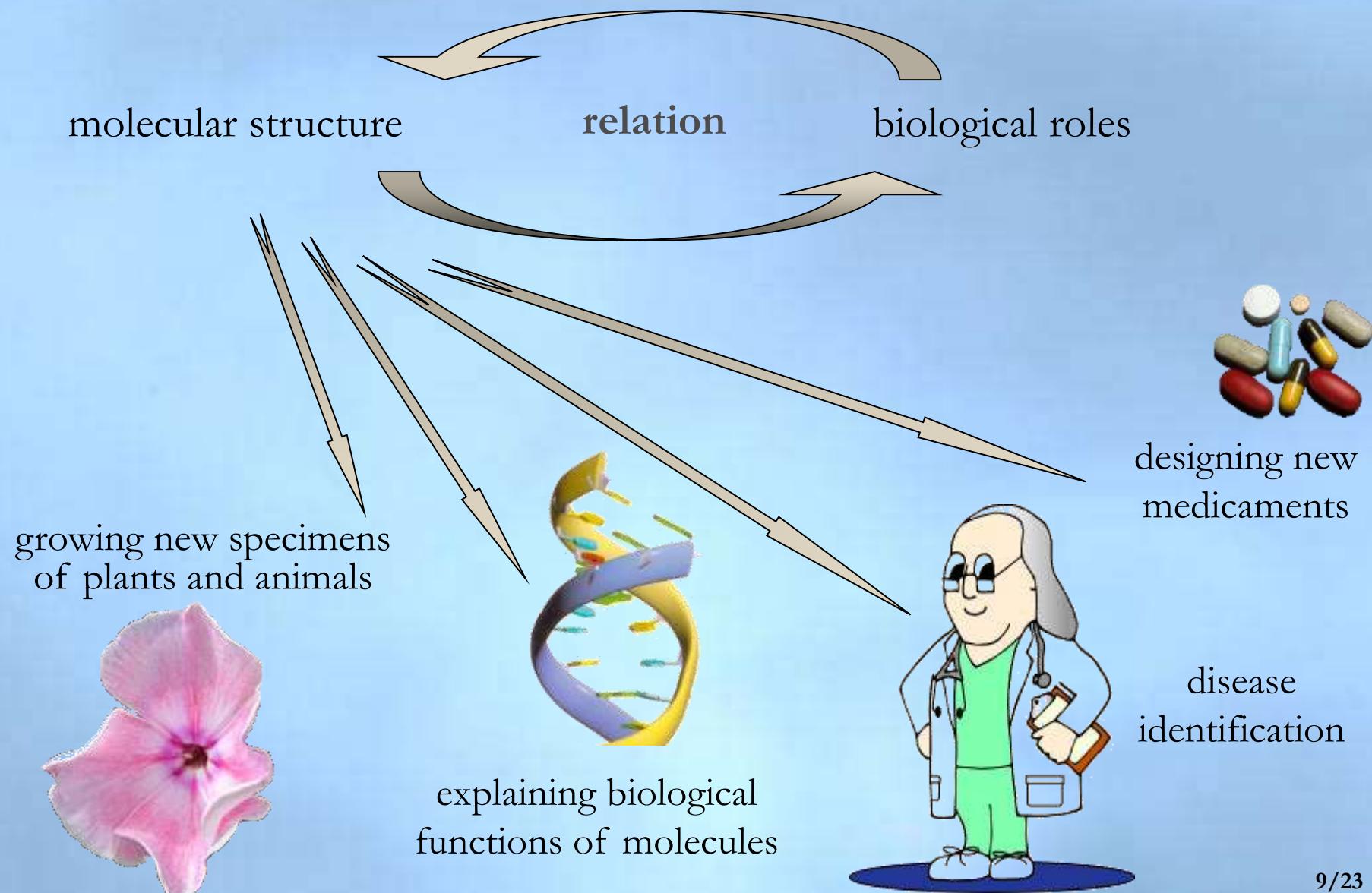


**Tertiary** structure is a 3D shape an entire chain assumes.

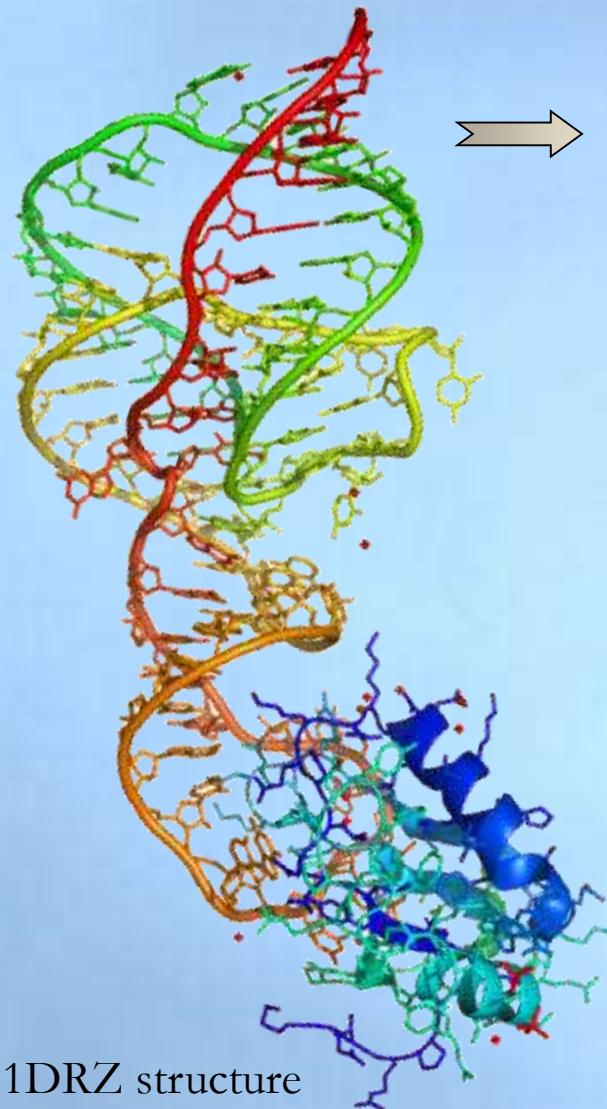
Determines molecule physiological roles inside a cell.

Describes ribose ring folding, torsion angles, base-pairing (canonical, Watson-Crick A:U, C:G, non-canonical>100 types).

# Goals of structure determination



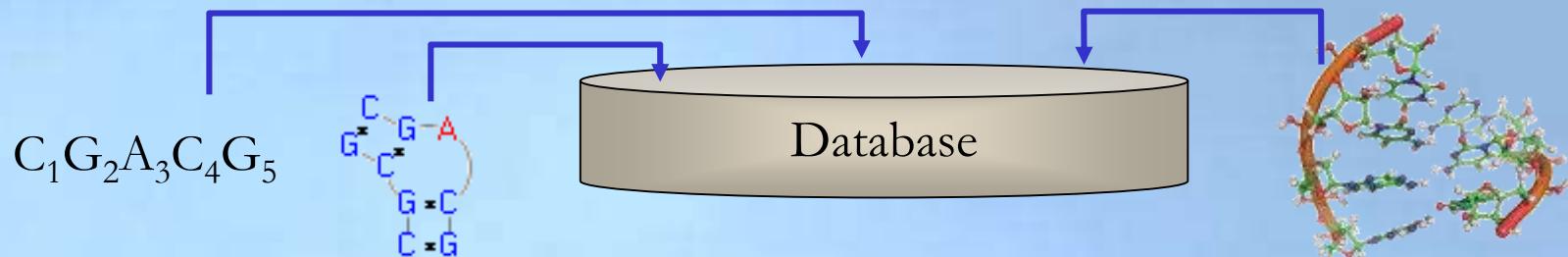
# What follows structure determination...



- Protein Data Bank (PDB)  
<http://www.rcsb.org/pdb/>
- Nucleic Acid Database (NDB)  
<http://ndbserver.rutgers.edu/>
- The RNA Structure Database (RNABase)  
<http://www.rnabase.org/>
- The European Ribosomal RNA Database (rRNA)  
<http://www.psb.ugent.be/rRNA/>
- A Database of Metal Ion Binding Sites in RNA (MERNA)  
<http://merna.lbl.gov/>
- The Structural Classification of RNA (SCOR)  
<http://scor.lbl.gov/scor.html>

# RNA fragments database: main concept

Construct a relational database storing the information about primary, secondary and tertiary structures of RNA molecules and / or their fragments.



Find an answer to the following problems / questions...

➡ look for 3D structures of RNA with the given sequence ...

**C G C G U A C ...**

➡ look for 3D structures of RNA with  
the given secondary structure motif ...



# RNA fragments database: main concept

Look for 3D structure of the RNA fragment defined by:

- (a) Sequence and / or
- (b) Secondary structure

1FFK 5S rRNA *H. Morismortui*

The diagram shows a complex, multi-stranded RNA molecule with various loops and junctions, colored in shades of green, blue, and yellow.

RNA FRABASE search

```
>strand1  
GGCCACAGCGGU  
((...(((  
>strand2  
ACCAGCGUUCGG  
))))...(((  
>strand3  
CCGGUUCGCC  
))))....))
```

three-way junction fragment

A simplified secondary structure diagram showing a three-way junction. The strands are labeled with numbers 10, 70, and 110 at their vertices. The strands are colored purple, green, and orange respectively, matching the colors used in the 3D structure above.

An alignment diagram showing the sequence of the RNA fragment. The sequence is color-coded to match the strands in the 3D structure. The sequence is:

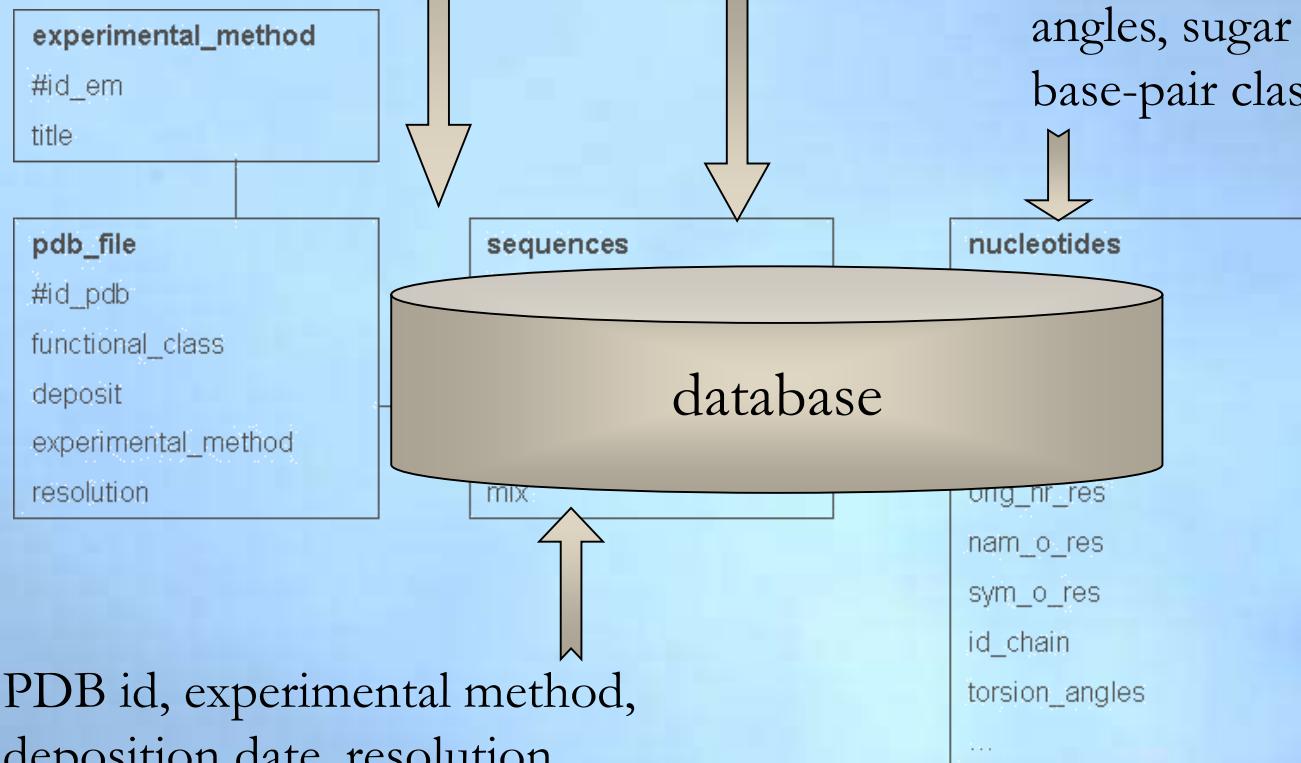
uUAGGC<sub>GG</sub>CCACAGCGGU<sub>GGG</sub>GUUGCCU<sub>CCC</sub>GUACCCAUC  
CCGAACACGGAAAGAUAA<sub>AG</sub>CCCACCAGCGUUC<sub>CGG</sub>GGAGUA  
CUGGAGUGCGAGCCUCUGGGAAAC<sub>CCGG</sub>UUCGCCGCCA  
CC  
... ((((((... (((((... (((((... (....  
..) ...)))...) ...) ....)))))))...) ... ((((...  
((((((.((...)))))))....)))...) )))...)))

# RNA fragments database: contents

Primary structure (sequence)  
given in IUPAC codes.

Secondary structure given  
in dot-bracket notation.

Atom coordinates, torsion  
angles, sugar pucker parameters,  
base-pair classification



PDB id, experimental method,  
deposition date, resolution.

# Primary structure encoding

IUPAC-IUB codes for RNA

Code	Description
A	Adenine
C	Cytosine
G	Guanine
U	Uracil
R	Purine {A, G}
Y	Pyrimidine {C, U}
M	{C, A}
K	{U, G}
W	{U, A}
S	{C, G}
B	{C, U, G} (not A)
D	{A, U, G} (not C)
H	{A, U, C} (not G)
V	{A, C, G} (not U)
N	any base {A, C, G, U}

The sequence in the database is defined with the code over 4-letter alphabet **{A, C, G, U}**.

Upper-case letter codes for unmodified residue, lower-case letter codes for modified residue.

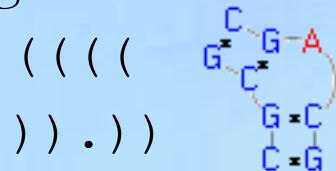
Pattern sequence can be defined with **IUPAC-IUB** code for RNA.

Search engine uses regular expression matching procedures to find the pattern within the deposited RNA structures.

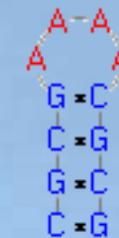
# Secondary structure encoding

Dot-bracket notation: . = unpaired residue ( ) = base-pair

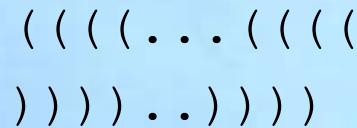
bulge



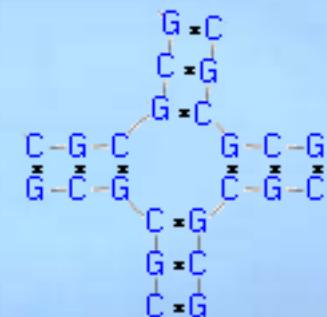
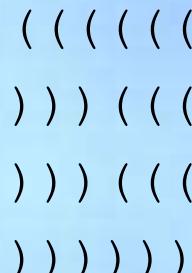
## hairpin loop



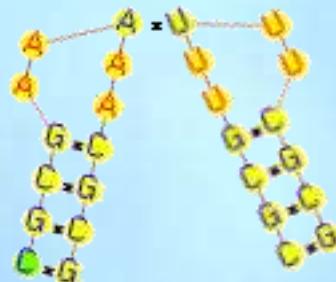
## internal loop



## 4-way junction



# kissing loops



# Filtering structures

Look for structures that suit the pattern and satisfy the given requirements:

- (a) have been determined by the given method (e.g. X-ray, NMR),
- (b) having torsion angle value(s) within the given range,
- (c) with the given sugar pucker parameters, etc.

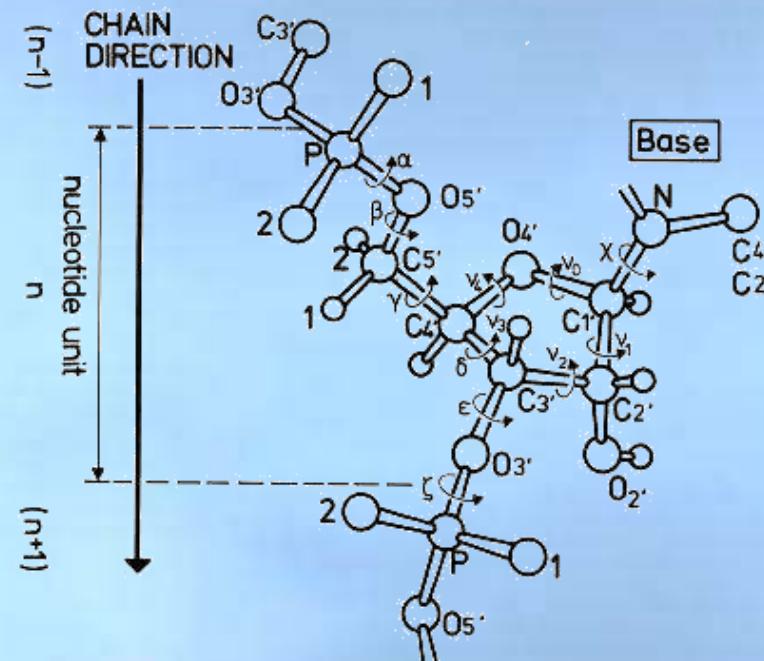
## example

Res. 2 contains torsion angles similar to those in B-DNA structure

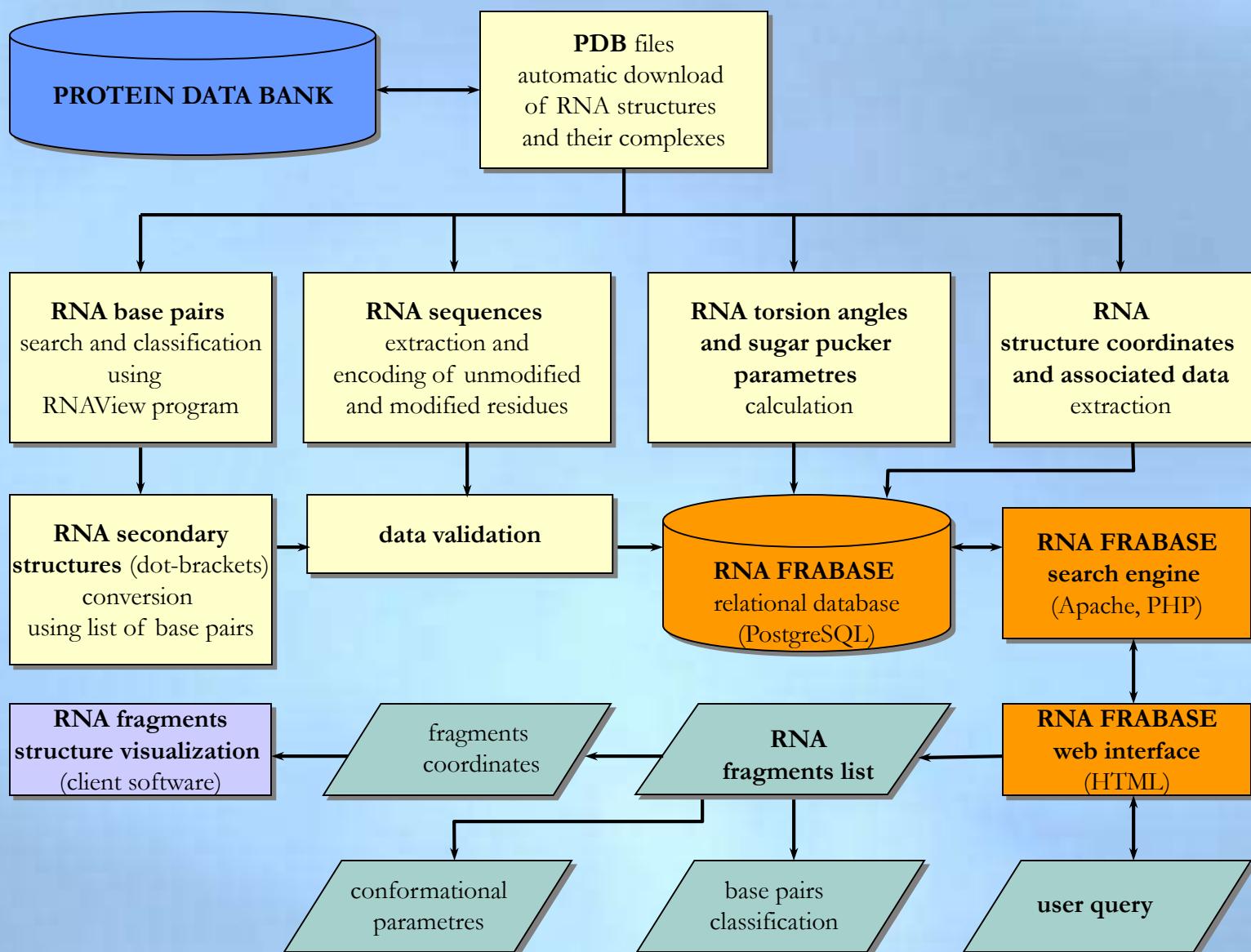
2	N ALPHA	-46	-36
2	N BETA	131	141
2	N GAMMA	33	43
2	N DELTA	134	144
2	N CHI	-137	-67
2	N EPSILON	-138	-128
2	N ZETA	-163	-157

Residue 3 is in anti & C3'-endo conformation. Res. 4 is in syn & C2'-endo conformation (conformations characteristic for left-handed structures).

3	C CHI	-180	-110
3	C P	0	36
4	G CHI	50	90
4	G P	144	180



# RNA fragments database: scheme



# RNA FRABASE version 1.0

## RNA FRAgments search engine & dataBASE

Search

About

Help

PDB list

References

Links

Contact us

You are 857. visitor.

Currently 1 user(s) on-line.

Supported by  
 Laboratory of Structural  
 Chemistry of Nucleic  
 Acids,  
 Institute of Bioorganic  
 Chemistry,  
 Polish Academy of  
 Sciences

**RNA FRABASE:** an engine with database to search the three-dimensional fragments within 3D RNA structures using as an input the sequence(s) and / or secondary structure(s) given in the dot-bracket notation.

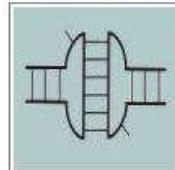
Database contains (updated **2007-07-30**):

- RNA sequences and secondary structures, described in dot-bracket notation, derived from the **1065** PDB-deposited RNA structures and their complexes,
- atom coordinates of the unmodified and modified PDB-deposited RNA structures,
- torsion angle values, sugar pucker parameters

### Sequence(s) and / or secondary structure(s):

Load example: 1 2 3 4 5 **6**

```
#Kissing loops
>strand1
(((.|||||)))
>strand2
((.|||||)))
```

**Search****Advanced Search**

Check for modified residues sensitive search:

Experimental method: X-Ray Diffraction   
 NMR   
 Electron Microscopy   
 Other

Resolution  $\leq$   Å

Limits for pseudorotation parameters, sugar pucker amplitude and torsion angles:

#all	alpha	-180	180
#all	beta	-180	180
#all	gamma	-180	180
#all	delta	-180	180
#all	epsilon	-180	180
#all	zeta	-180	180
#all	chi	-180	180
#all	p	-180	180
#all	vmax	-180	180
#all	v0	-180	180
#all	v1	-180	180

**Search****Reset**

## Submission results

Number of matching fragments: 33

No.	PDB id	sequence	secondary structure	chain	start	end	method	class	date	A
<input type="checkbox"/>	1	1FFK	CAGGCAUCGACI CUGCUUGAUGC	G	((...[[[[[					
<input checked="" type="checkbox"/>	2	1JJ2	CAGGCAUCGACI CUGCUUGAUGC	G	((...[[[[[					
<input type="checkbox"/>	3	1K73	CAGGCAUCGACI CUGCUUGAUGC	G	((...[[[[[					

Base pairs classification for selected fragment(s) → Save table

residue I				residue II				classification		
PDI	name	symbol	chain	res no.	PDB name	symbol	chain	res no.	Westhof's notation	Saenger's notation
	C	C	0	414	G	G	0	426	+/- cis	XIX
									-/- cis	XX
									+/- cis	XIX
									+/- cis	XIX
									-/- cis	XX
									+/- cis	XX
									-/- cis	XX

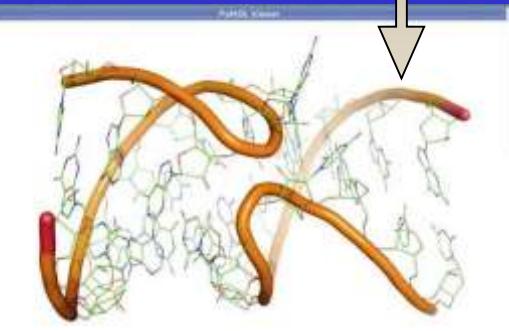
Atom coordinates for selected fragment(s) → Save table

No.	PDB id	sequence	secondary structure	chain	start	end	method	class	date	A
2	1JJ2	CAGGCAUCGACUG CUGCUUGAUGC	((...[[[[[)))	0	414	426	X-Ray	RIBOSOME	2001-07-03	2.4

Torsion angles for selected fragment(s) → Save table

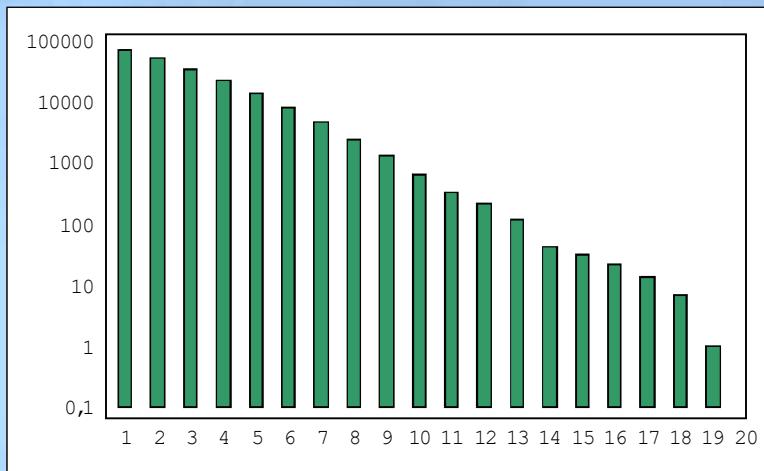
No.	PDB id	sequence	secondary structure	chain	start	end	method	class	date	A
2	1JJ2	CAGGCAUCGACUG CUGCUUGAUGC	((...[[[[[)))	0	414	426	X-Ray Diffraction	RIBOSOME	2001-07-03	2.4

name	symbol	str.	chain	res no.	α	β	γ	δ	ε	ζ	X	P	vmax	v0	v1	v2	v3	v4
C	C	(	0	414	-66.7	166.4	58.2	81.0	-152.0	-72.3	-162.3	11.7	42.6	5.2	-29.8	41.7	-39.7	21.7
A	A	(	0	415	-65.6	176.7	60.2	87.2	-157.2	-68.2	-166.2	14.5	36.5	2.7	-24.3	35.4	-34.9	20.3
G	G	(	0	416	-68.2	-177.4	50.5	85.6	-126.4	72.8	-148.1	15.7	35.2	2.0	-22.9	33.8	-34.0	20.2
G	G	.	0	417	-127.4	-100.0	76.8	130.9	-93.1	176.0	-93.9	151.9	31.0	-23.2	31.3	-27.4	14.7	5.2
C	C	[	0	418	-54.3	-142.6	57.6	82.8	-159.8	-70.2	-165.7	12.0	38.9	4.7	-27.2	38.0	-36.6	20.1
A	A	[	0	419	-71.5	-172.2	50.9	85.4	-149.0	-74.7	-154.8	9.7	38.5	6.1	-28.0	38.0	-35.6	18.5



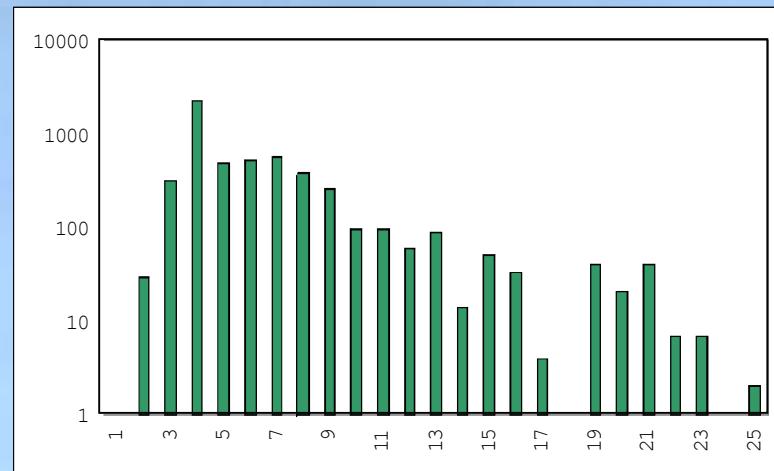
# Selected query statistics

A



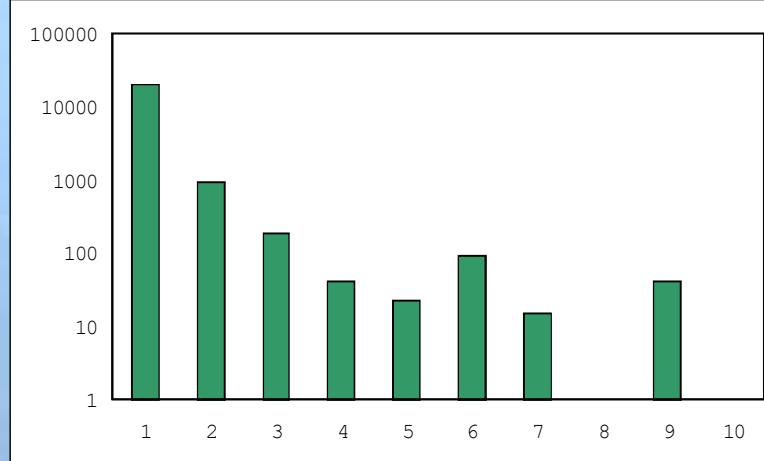
Number of RNA duplexes for different number of base-pairs.

B



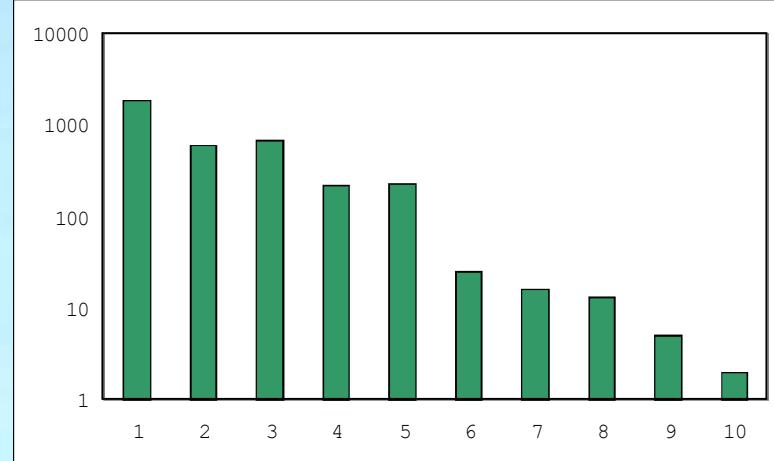
Number of internal loops for different number of unpaired residues.

C



Number of bulges for different number of unpaired residues.

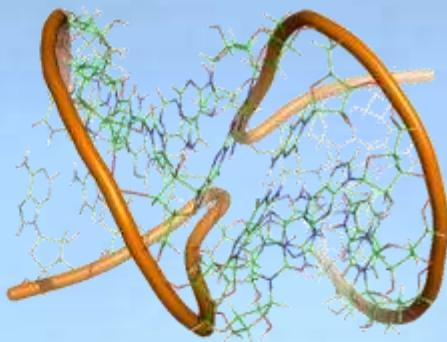
D



Number of symmetric internal loops for different number of unpaired residues.

# Kissing loop search: selected results

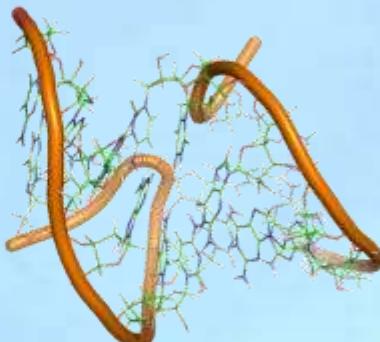
>Strand1  
(([[[[[[[ ] )  
>Strand2  
((( ] ] ] ] ] ])))  
  
PDB ID : 1BJ2



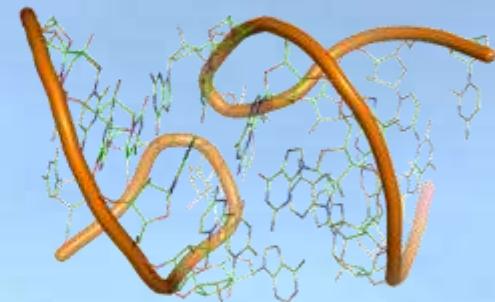
>Strand1  
(((..[ [ ))  
>Strand2  
(((..] ] ))))  
  
PDB ID: 1F5U



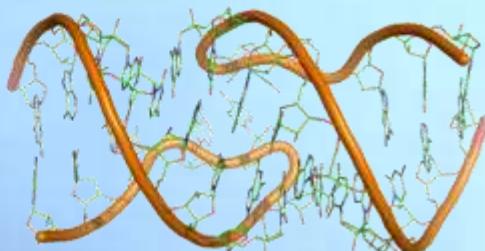
>Strand1  
(.....[ )  
>Strand2  
(.].....)  
  
PDB ID: 1KIS



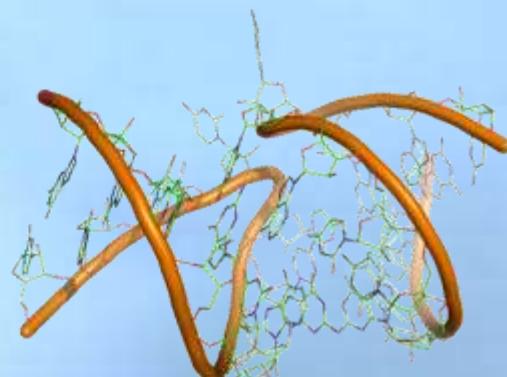
>Strand1  
((..[ [[ [ ))  
>Strand2  
((..] ] ] ..))  
  
PDB ID: 1JZX



>Strand1  
(((..[ [[[[[ [ .)))  
>Strand2  
(((..] ] ] ] ] .)))  
  
PDB ID: 1ZCI



>Strand1  
(((..[ [[[[[ [ ))  
>Strand2  
(((..] ] ] ] ] .)))  
  
PDB ID: 1VQO



# Selected search results

Chemical structures and sequence contexts for selected search results																									
I		II		III		IV		V		VI		VII		VIII		IX		X		XI		XII			
A 5'-A C-3' 3'-U-G-5'	8 ⑧	C 5'-U G-3' 3'-A-C-5'	2 ②	UGU 3'-G C-5'	0	UU 5'-U A-3'	0	GAC 3'-G G-5'	0	U 5'-C G-3'	0	AA 3'-U U-5'	0	AA-AA CA	AAUCUC	5'-G-C-3' 3'-C G-5'	0	A 5'-G G-3' 3'-C-C-5'	84 ⑧	AGC 5'-G G-3' 3'-U---C-5'	0	G-G-A AACAA	5'-G U-3' 3'-C A-5'	0	
A 5'-R Y-3' 3'-Y-R-5'	192 56	C 5'-Y R-3' 3'-R-Y-5'	66 56	UGU UUC	0	UU 3'-R R-5'	0	GAC 3'-R Y-5'	67 67	U 5'-Y R-3'	0	AA 3'-Y Y-5'	0	AA-AA CA	AUCUC	5'-R-Y-3' 3'-Y R-5'	0	A 5'-R R-3' 3'-Y-Y-5'	144 1	AGC 5'-R R-3' 3'-Y---Y-5'	0	G-G-A AACAA	5'-R Y-3' 3'-Y R-5'	0	
A 5'-N N-3' 3'-N-N-5'	848 323	C 5'-N N-3' 3'-N-N-5'	0 323	UGU UUC	0	UU 3'-N N-5'	0	GAC 3'-N N-5'	0	U 5'-N N-3'	0	AA 3'-N N-5'	148 39	AA-AA CA	AUCUC	5'-N-N-3' 3'-N N-5'	0	A 5'-N N-3' 3'-N-N-5'	150 848	AGC 5'-N N-3' 3'-N---N-5'	0	G-G-A AACAA	5'-N N-3' 3'-N N-5'	0	
R 5'-A C-3' 3'-U-G-5'	8 8	Y 5'-U G-3' 3'-A-C-5'	4 4	YRY 3'-G G-5'	0 0	YY 3'-A U-5'	0	RRY 3'-G G-5'	0	Y 5'-U C-3'	0	C 3'-G U-5'	31 31	RR 3'-U U-5'	0	RR-RR YRYYY	5'-G-C-3' 3'-C G-5'	0	R 5'-G G-3' 3'-C-C-5'	84 2	RRY 5'-G G-3' 3'-U---C-5'	0	R-R-R RRYRR	5'-G U-3' 3'-C A-5'	0
R 5'-R Y-3' 3'-Y-R-5'	258 82	Y 5'-Y R-3' 3'-R-Y-5'	82 82	YRY YYY	0	YY 3'-R Y-5'	0	RRY 3'-R R-5'	43 43	Y 5'-Y Y-3'	0	R 3'-R Y-5'	140 0	RR 3'-Y Y-5'	0	RR-RR RYYYY	5'-R-Y-3' 3'-Y R-5'	0	R 5'-R R-3' 3'-Y-Y-5'	203 44	RRY 5'-R R-3' 3'-Y---Y-5'	0	R-R-R RRYRR	5'-R R-3' 3'-R R-5'	0
R 5'-N N-3' 3'-N-N-5'	1206 686	Y 5'-N N-3' 3'-N-N-5'	0 39	YRY YYY	0	YY 3'-N N-5'	0	RRY 3'-N N-5'	73 73	Y 5'-N N-3'	0	RR 3'-N N-5'	73 73	RR-RR RYYYY	5'-N-N-3' 3'-N N-5'	0	R 5'-N N-3' 3'-N-N-5'	370 1206	RRY 5'-N N-3' 3'-N---N-5'	19 19	R-R-R RRYRR	5'-N N-3' 3'-N N-5'	0		
N 5'-A C-3' 3'-U-G-5'	8 8	N 5'-U G-3' 3'-A-C-5'	44 44	NNN 3'-G C-5'	108 108	NN 3'-A U-5'	9	NNN 3'-G G-5'	104 104	N 5'-U C-3'	0	NN 3'-G U-5'	63 63	NN 3'-U U-5'	5 5	NN-NN NNNNNN	5'-G-C-3' 3'-C G-5'	0	N 5'-G G-3' 3'-C-C-5'	274 14	NNN 3'-U---C-5'	0	N-N-N NNNNNN	5'-G U-3' 3'-C A-5'	0
N 5'-R Y-3' 3'-Y-R-5'	357 529	N 5'-Y R-3' 3'-R-Y-5'	279 157	NNN 3'-Y Y-5'	279 157	NN 3'-R R-5'	350 350	NNN 3'-R Y-5'	579 579	NN 5'-Y Y-3'	0	NN 3'-Y Y-5'	310 310	NN-NN NNNNNN	5'-R-Y-3' 3'-Y R-5'	34 265	N 5'-R R-3' 3'-Y-Y-5'	464 98	NNN 3'-Y---Y-5'	98 NNNNNN	N-N-N NNNNNN	5'-R Y-3' 3'-Y R-5'	12		
N 5'-N N-3' 3'-N-N-5'	1883 1883	N 5'-N N-3' 3'-N-N-5'	0 667	NNN 3'-N N-5'	667 753	NN 3'-N N-5'	0	NNN 3'-N N-5'	1819 597	N 5'-N N-3'	0	NN 3'-N N-5'	597 597	NN-NN NNNNNN	5'-N-N-3' 3'-N N-5'	0	N 5'-N N-3' 3'-N-N-5'	947 1883	NNN 3'-N---N-5'	188 111	N-N-N NNNNNN	5'-N N-3' 3'-N N-5'	111		

## References

---

M. Popenda, M. Błażewicz, M. Szachniuk, R.W. Adamiak (2007) RNA FRABASE version 1.0: an engine with a database to search for the three-dimensional fragments within RNA structures, *Nucleic Acids Research*, doi: 10.1093/nar/gkm786.

M. Błażewicz, RNA fragments search engine – project and implementation, *M.Sc. thesis* to be completed in 2008

The database is freely accessible at:

<http://rnafrabase.ibch.poznan.pl>