# Progress toward Predicting Viral RNA Structure from Sequence:

# How Parallel Computing can Help Solve the RNA Folding Problem
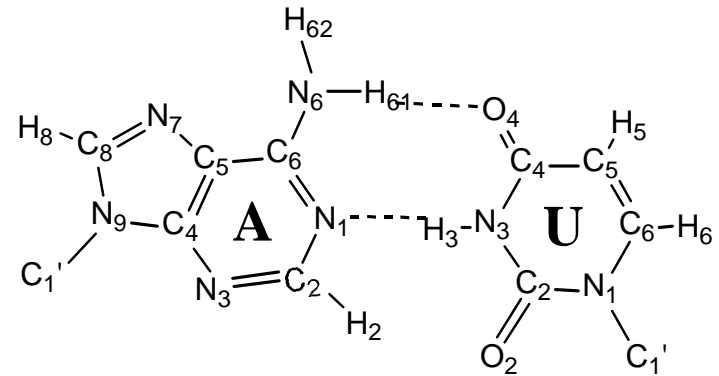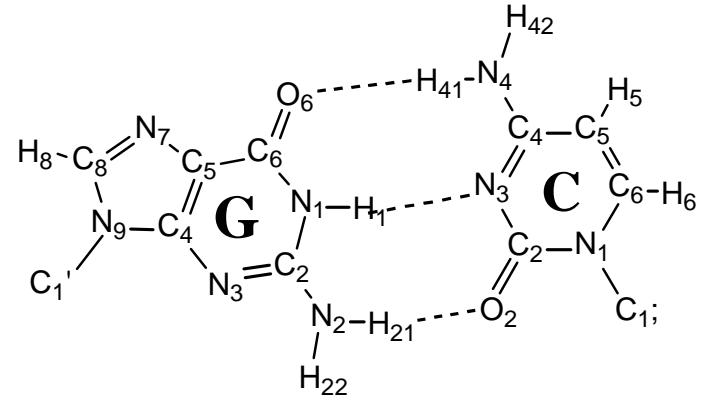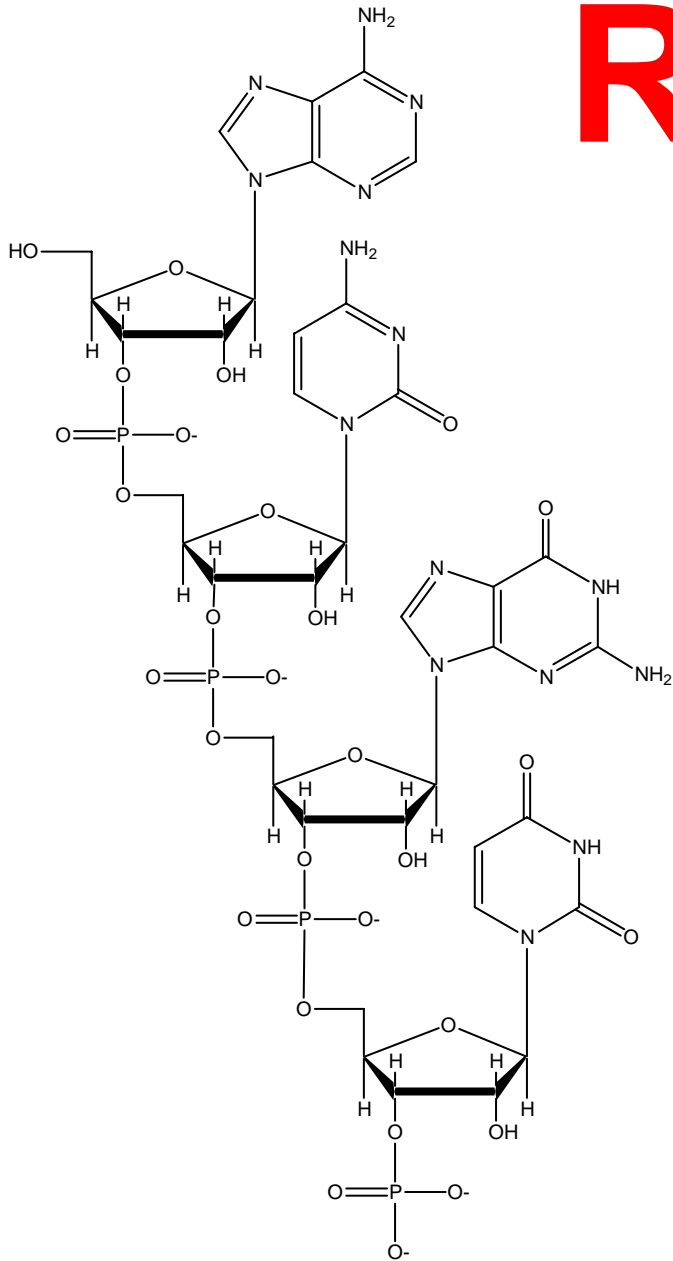
**Susan J. Schroeder**
**University of Oklahoma**
**October 7, 2008**

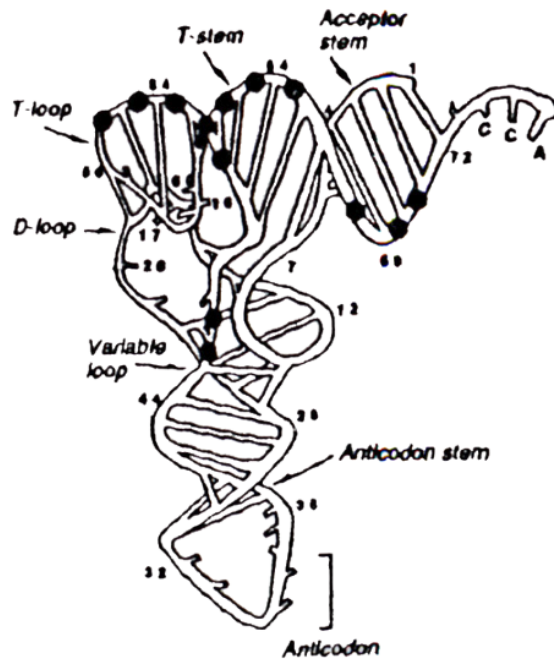We finished the genome map,
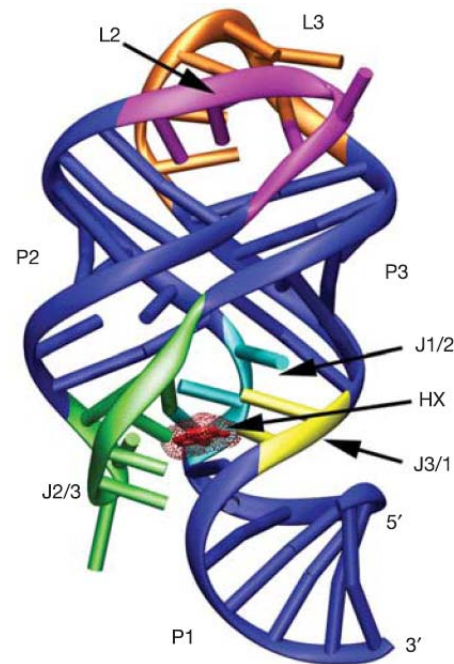now we can't figure out how to fold it!

Science (1989) **243,** p.786

# RNA

# Sequence → Structure → Function

5'GCGGAUUUAG$^{2M}$
CUCAGU$^{DH}$U$^{DH}$GGG
AGAGCG$^{M2}$CCAGA
C$^{OM}$UG$^{OM}$AAG$^{Y}$AU$^{PS}$
C$^{5M}$UGGAGG$^{7M}$UC
C$^{5M}$UGUGU$^{5M}$U$^{PS}$C
GA$^{1M}$UCCACAGAA

UUCGACCA

5'GGACAUAUAAU
CGCGUGGAUAUG
GCACGCAAGUUU
CUACCGGGCACC
GUAAAUGUCCGA
CUAUGUCCA



**tRNA**



**Guanine Riboswitch**

Batey, R. et al, 2004 Nature vol. 432, p. 412

# RNA Folding Problem



- **Folding a polymer with negative charge**
- **Watson - Crick base pairing**
- **Hierarchical folding**  **1° ➡ 2° ➡ 3°**

**1°** ➡ **2°** ➡ **3°**
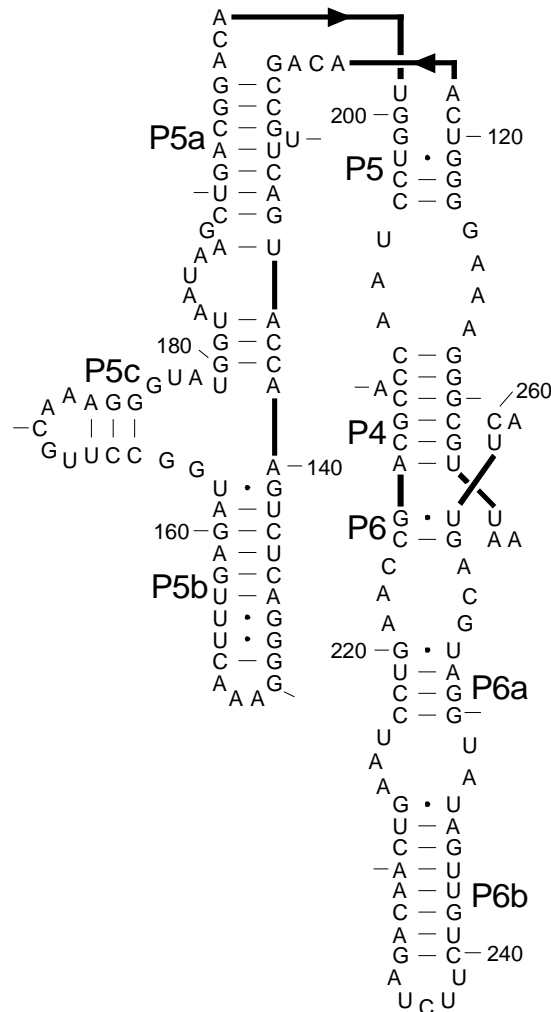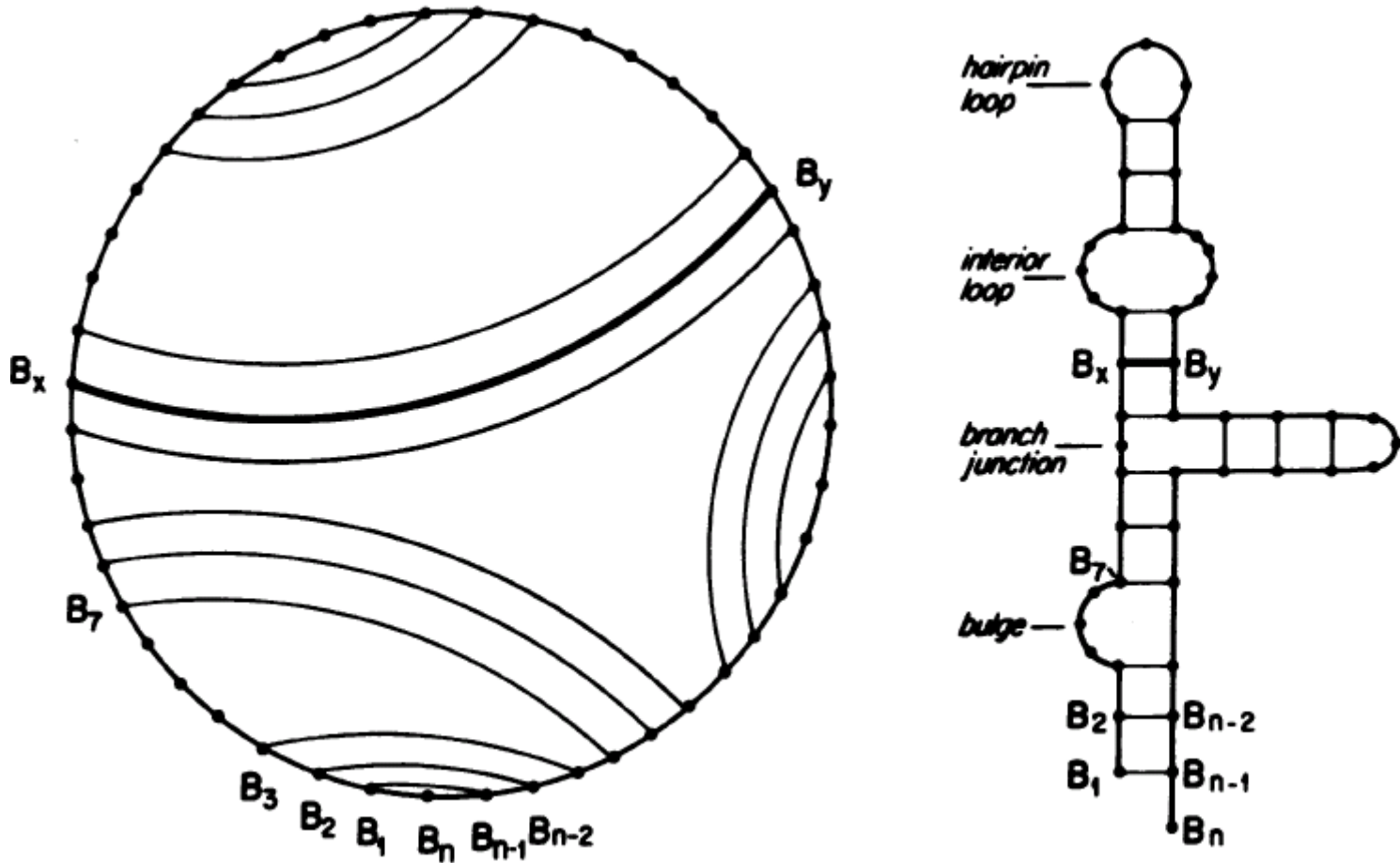
AAUUGCGGGAAAGGGGUCAA
CAGCCGUUCAGUACCAAGUC
UCAGGGGAAACUUUGAGAUG
GCCUUGCAAAGGGUAUGGUA
AUAAGCUGACGGACAUGGUC
CUAACCACGCAGCCAAGUCC
UAAGUCAACAGAUCUUCUGU
UGAUAUGGAUGCAGUUCA

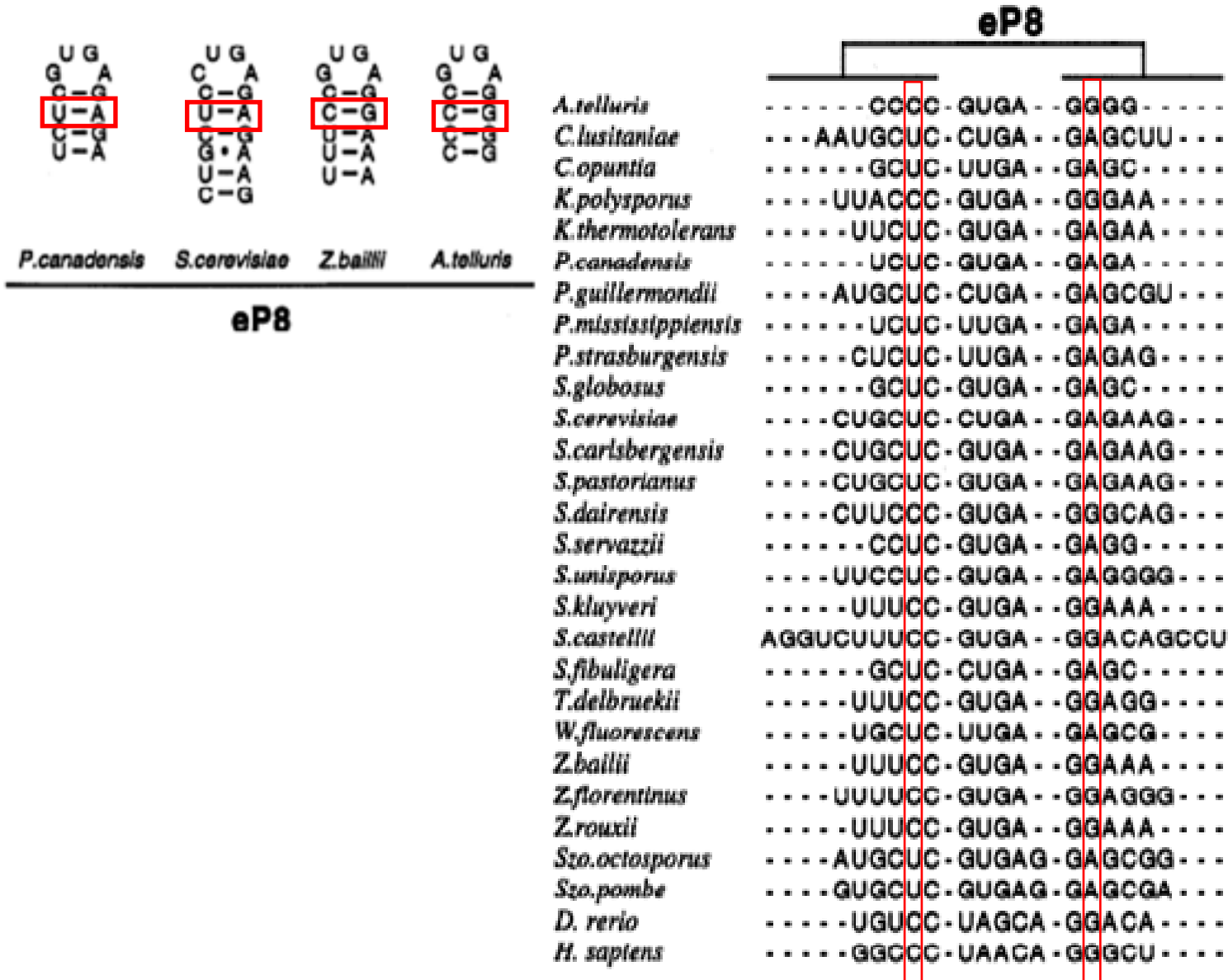Waring &Davies 1984
Gene 28:277

Cate et al. 1996
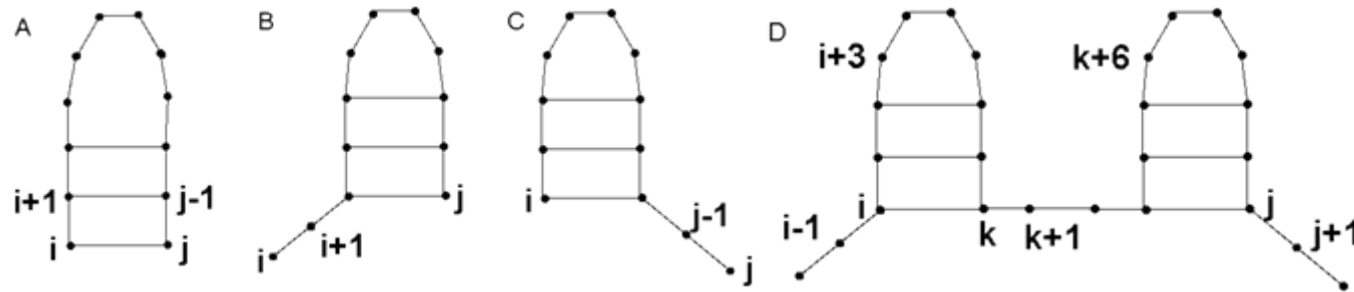Science 273:1678

# RNA Secondary Structure has Helices and Loops



Nussinov & Jacobson 1980 PNAS v 11 p 6310 Fig.1

# An example of phylogenetic alignment and structure prediction for RNAse P
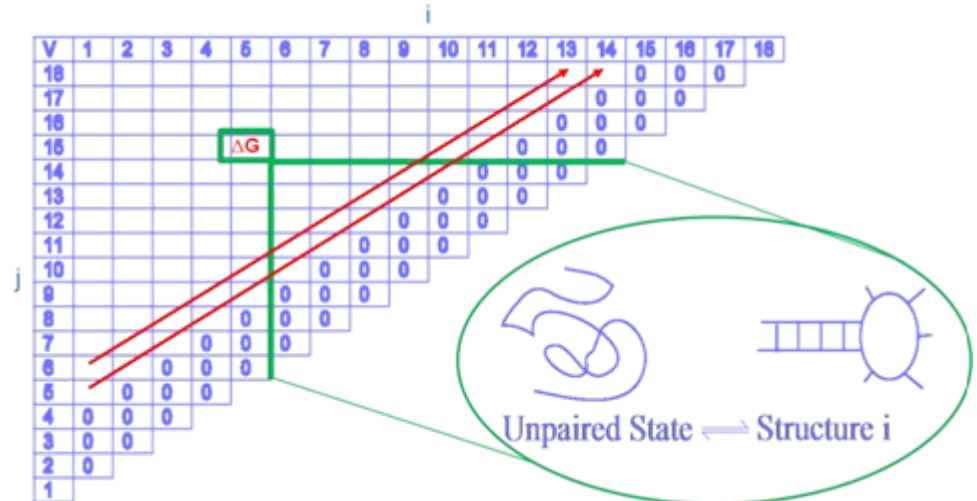
# How Dynamic Programming Algorithms Calculate RNA Secondary Structure

- **Stochastic context-free grammar defines possible base pairs as 1 of 4 possible cases**



- **Recursion statement finds maximum value for each small subset of RNA sequence**

- **Fill an array with scores for each substructure**

- **Traceback through the array to find the lowest free energy structure**



- **O(N²) memory storage**

- **O(N³) runtime**

# Websites for Folding Algorithms to Predict RNA 2°

**MFOLD**
http:// www.bioinfo.rpi.edu/applications/mfold

**Vienna package** http:// www.tbi/univie.ac.at/~ivo/RNA

**RNAStructure**  http://rna.urmc.rochester.edu

**SFOLD**  http://sfold.wadsworth.org

**PKNOTS**  http:// selab.wustl.edu

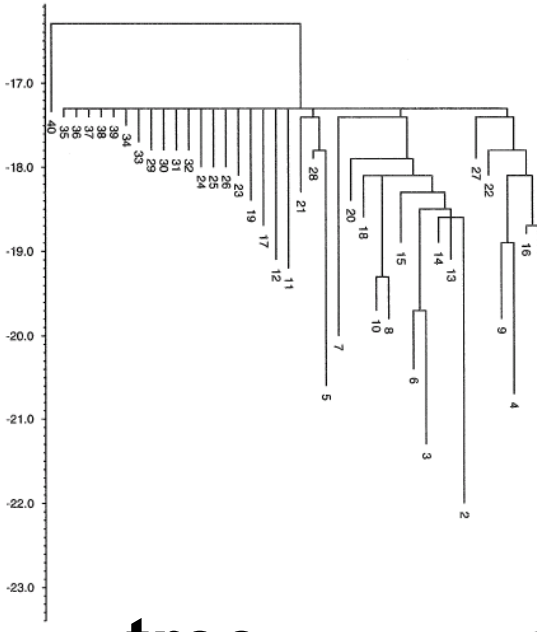**STAR4.4**
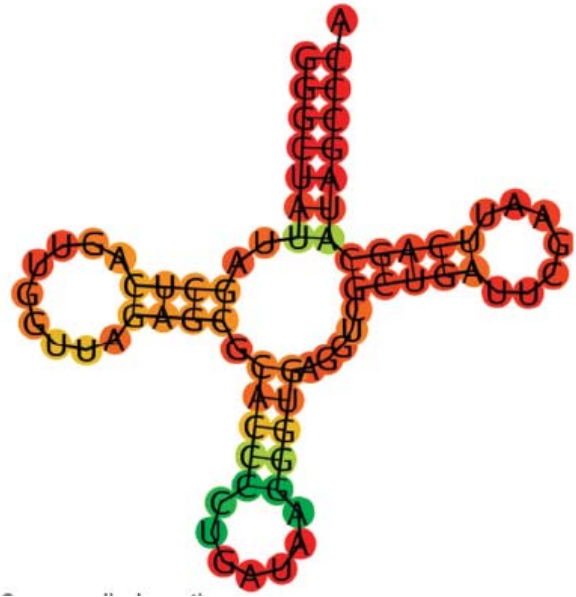http://biology.leidenuniv.nl/~batenburg/STRAbout.html

# Representations of RNA Secondary Structure

GGGCUAUUAGCUCAGUUGGUUAGAGCGCACCCCUGAUAAGGGUGAGGUCGCUGAUUCGAAUUCAGCAUAGCCCA
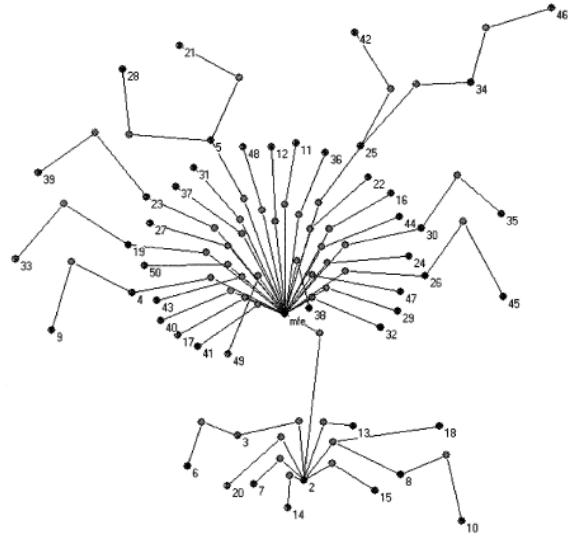(((((((..((((........)))).(((((.......)))))....(((((.......))))))))))).

## dots & parentheses
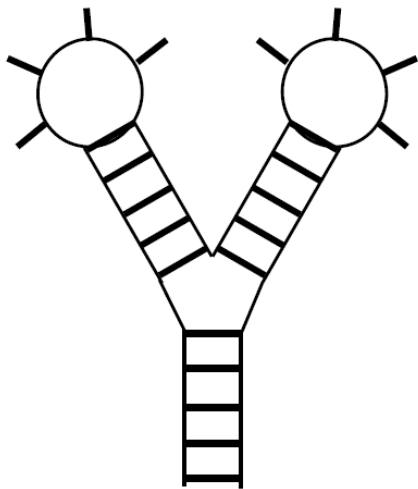


**tree**

**graph**

**merged landscape**

Wuchty 2003, Nucl. Acids Res. v 31, p 1115 Fig. 7
Gruber et al. 2008, Nucl. Acids Res. v 36, p. W73 Fig. 1
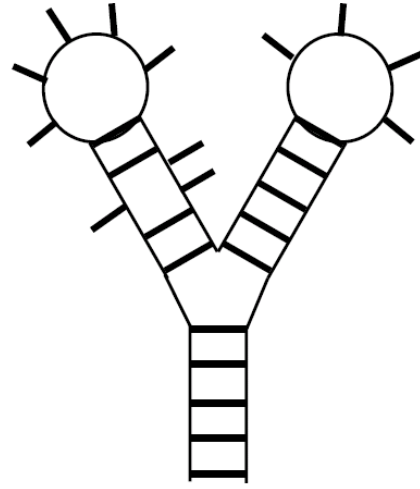
# RNAStructure Predicts Secondary Structure Well

| RNA | Lowest $\Delta G^o$ Structure | Best Suboptimal Structure |
|---|---|---|
| average | 73% | 87% |
| Group II introns | 88% | 94% |
| tRNA | 87% | 97% |
| 5 S rRNA | 74% | 96% |
| Group I introns | 69% | 84% |
| SRP RNA | 66% | 88% |
| Rnase P | 63% | 76% |
| 23 S rRNA | 55% | 61% |
| (as domains) | (74%) | (88%) |
| 16 S rRNA | 44% | 54% |
| (as domains) | (61%) | (76%) |

Mathews et al. (2004)  *Proc*. *Natl*. *Acad*. *Sci*. vol. 101, pp. 7287-7292
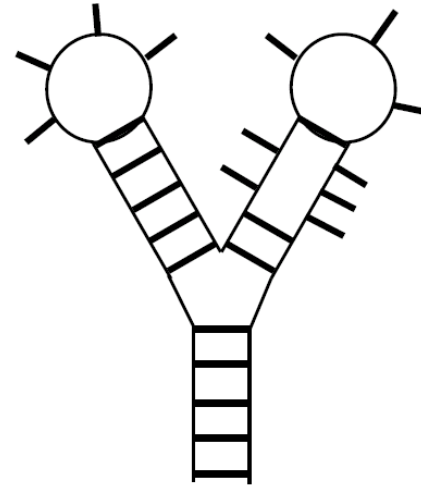
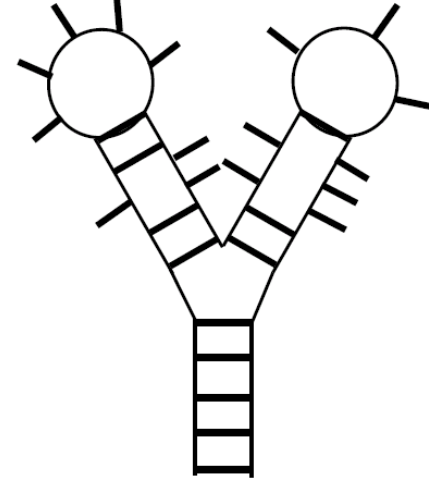# Mfold and  RNAStructure Sample Suboptimal Structures



Lowest Free Energy Structure
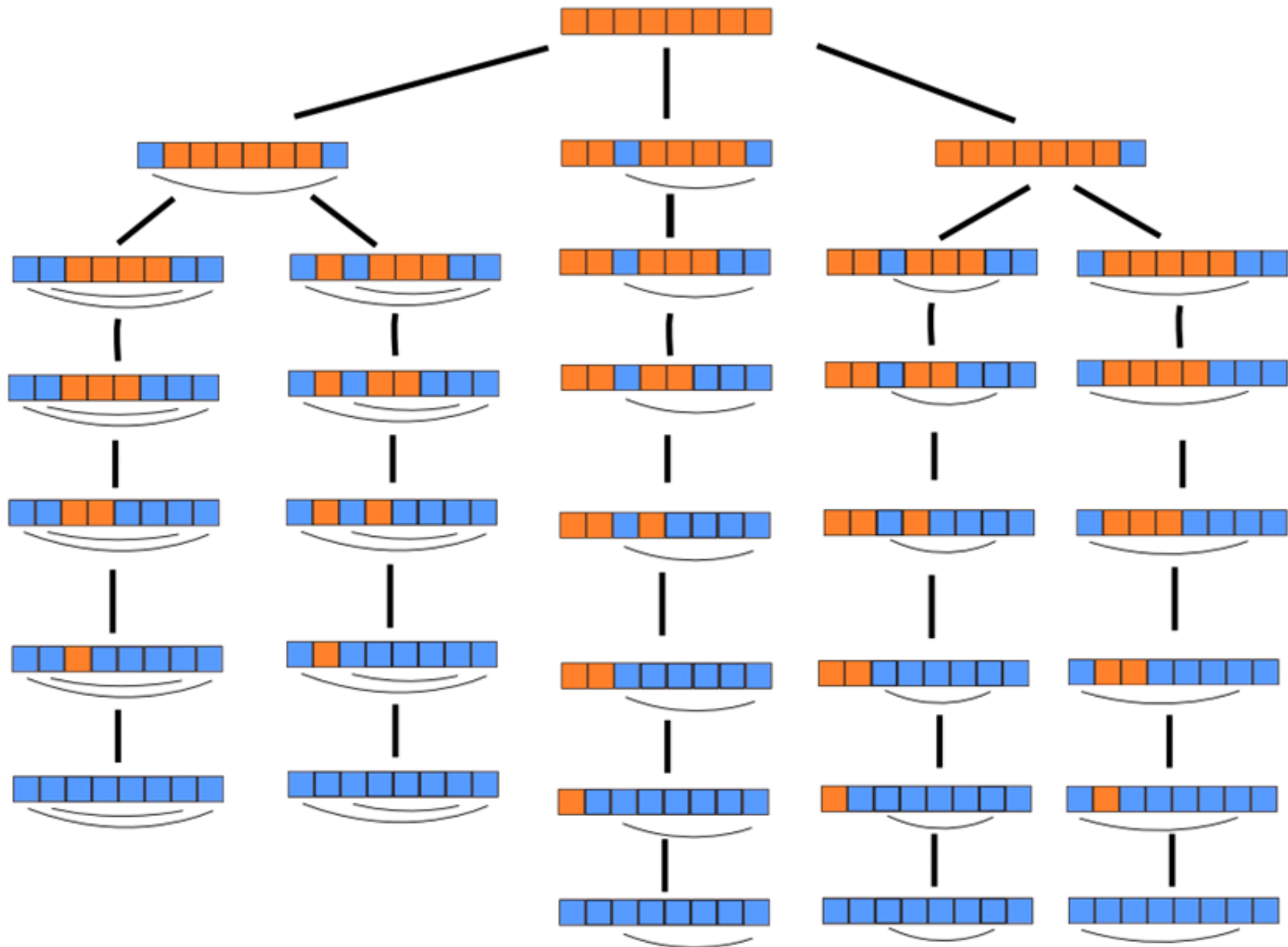
Suboptimal Structure 1

Suboptimal Structure 2

Structure that Cannot be Predicted

# How  Wuchty's Algorithm is like a Tree
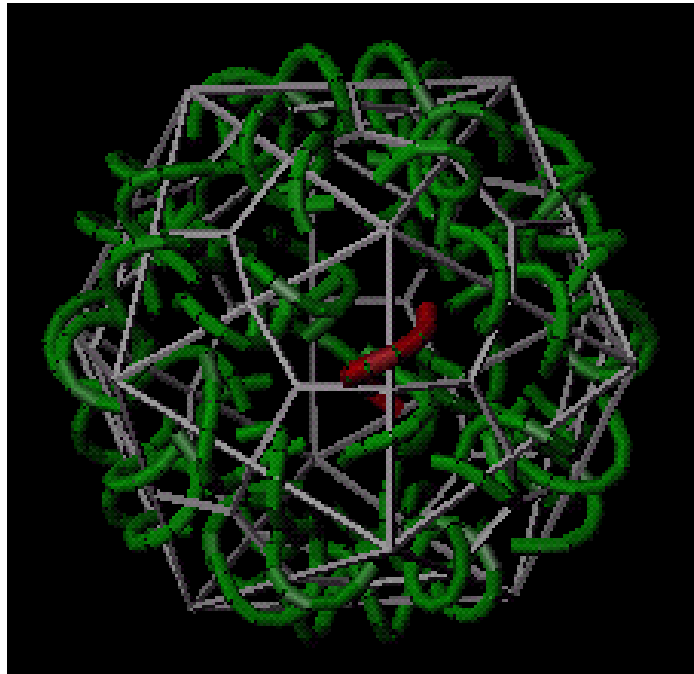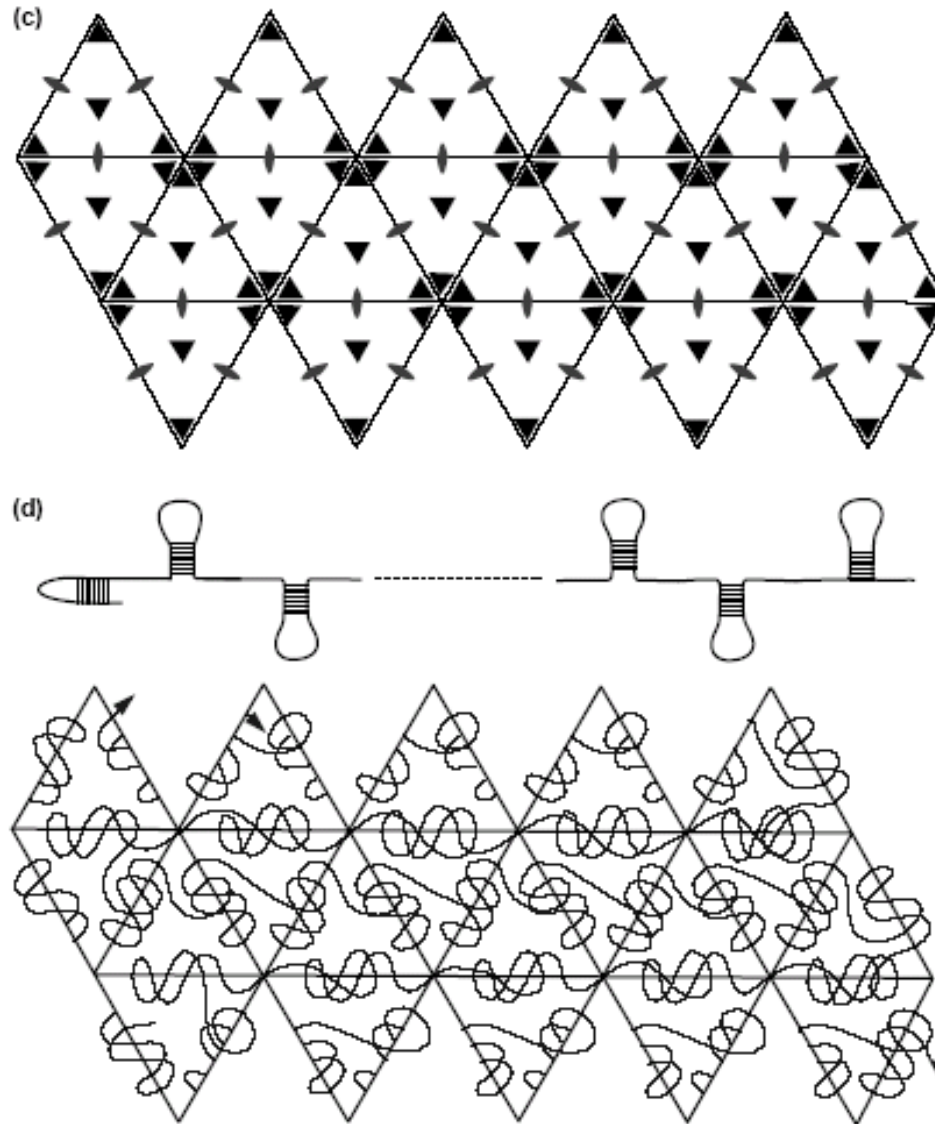
# STMV RNA folding problem



Figure reproduced from VIPER website
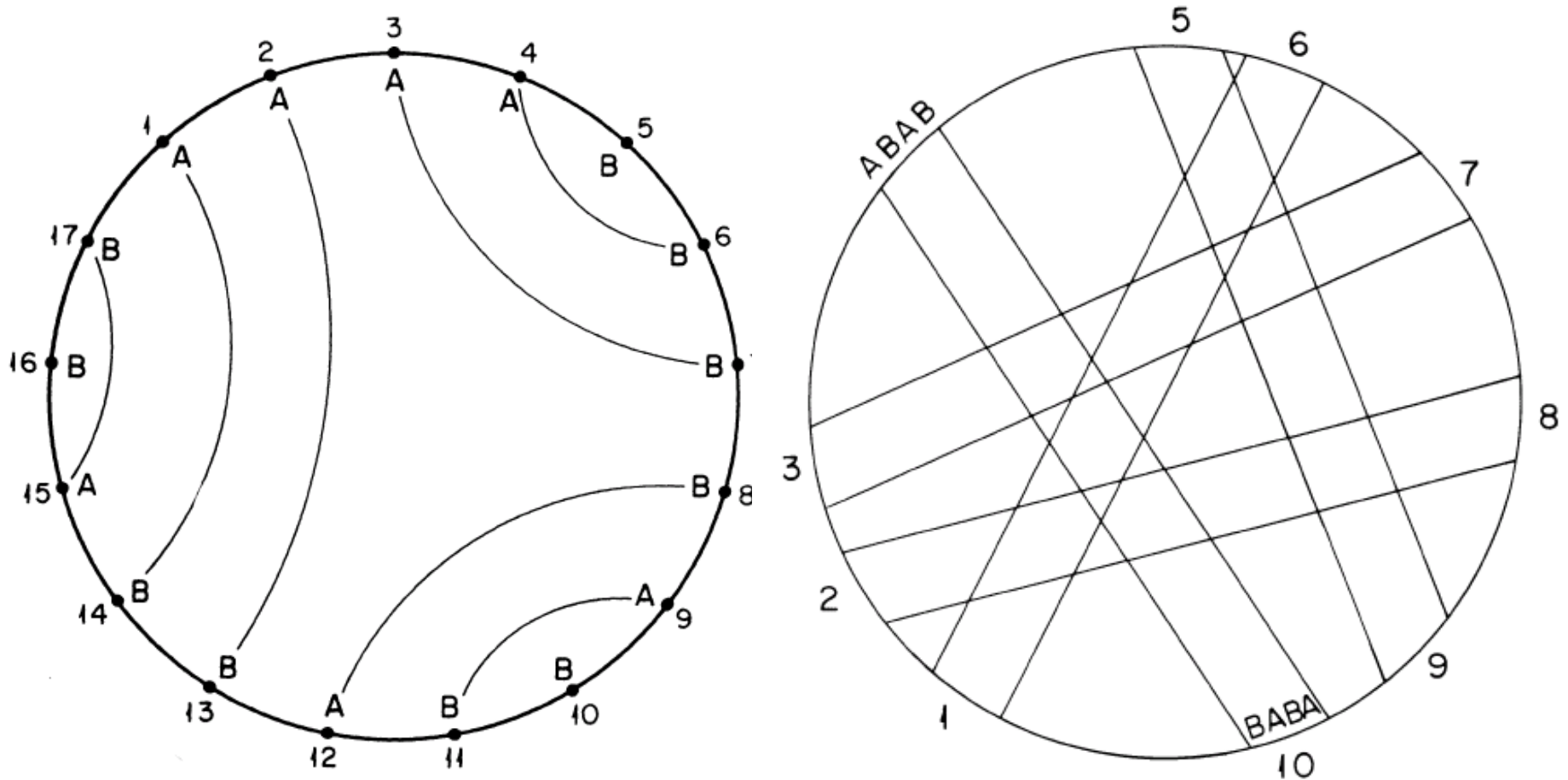Reddy et al. 2001

- **Crystal structure to 1.8 Å resolution (Larson et al., 1998)**
- **59% of the 1,058 RNA nucleotides are in helices**
- **RNA is icosahedrally averaged**
- **Identity of nucleotides in helices remains obscure**
- **Structure of 41% of the RNA remains unknown**
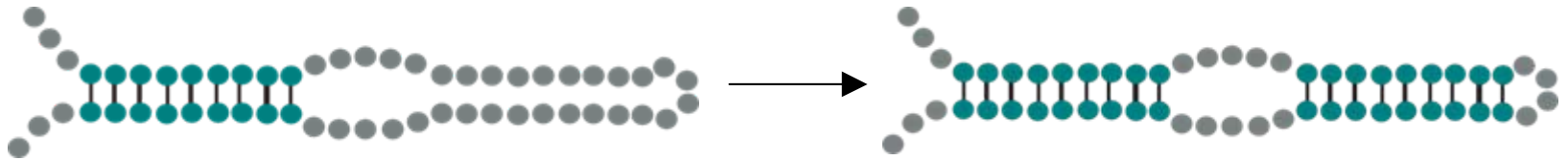
# Current model for STMV RNA



Larson, S. & McPherson, A. Current Opinions in Structural Biology, vol. 11, p. 61.

# Nussinov algorithms for maximizing matches and blocks



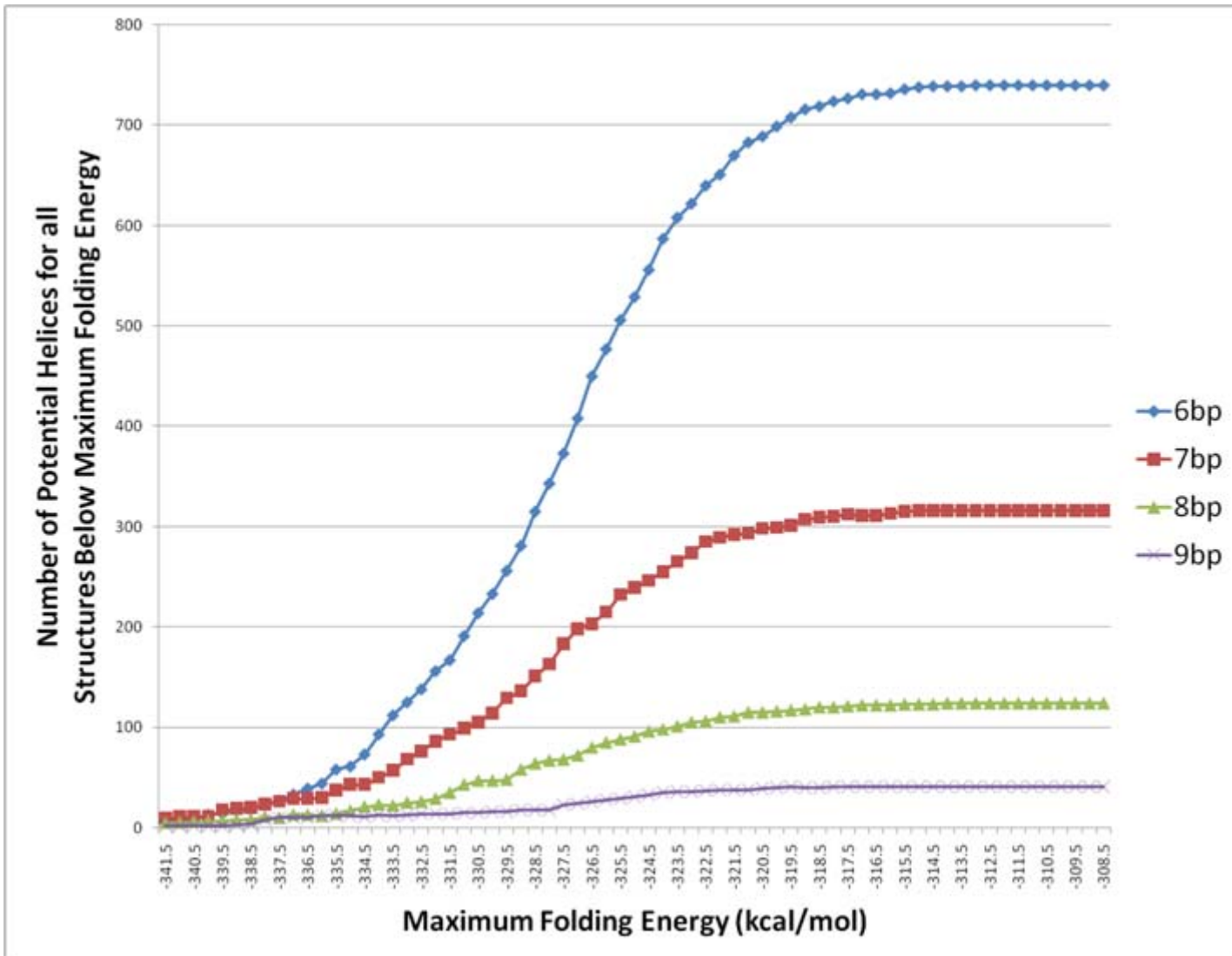Nussinov et al. 1978 SIAM v35, p. 71,78 Fig. 1, 5

# Combinatorial Search of STMV RNA

- **Locate potential helical structures between pairing bases *i* and *j***

- **Assemble non-overlapping potential helices *(i,j)* with *(p>j,q)* or *(k>i+l, k+l<q<j)***

- **Nested searching identifies**

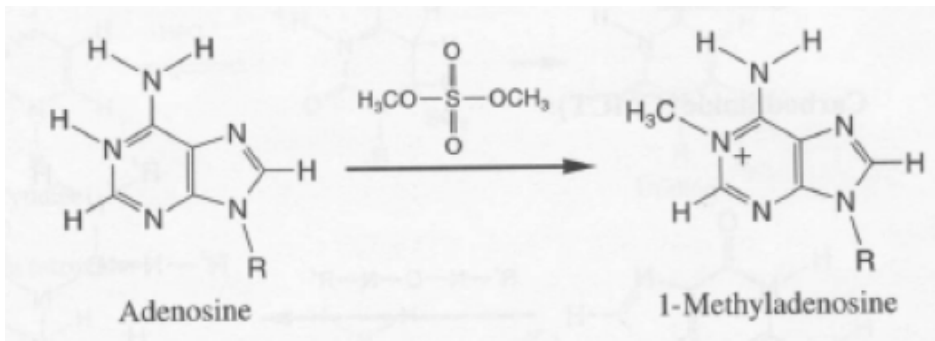  **"helices within helices"**



- **Over 144,000 perfect 6-pair helices, but no possible simultaneous combination of 30 helices in STMV RNA**

# Many more possible structures contain 30 imperfect helices in the STMV sequence

# Chemical modification data restrains possible base pairing



5'**<span style="color:red">A</span>**AGAUGU<sup>AAA</sup>CCAGGA3'
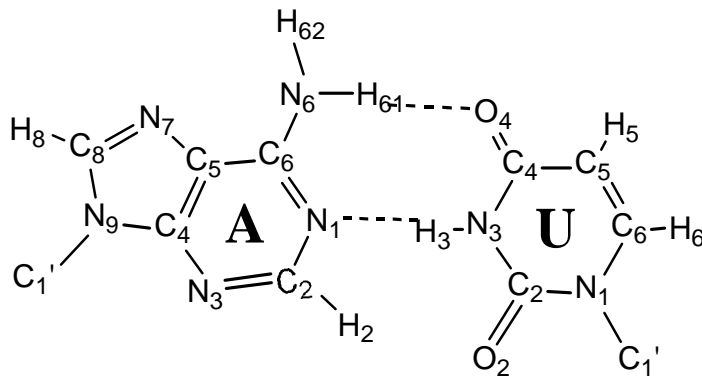3'CUGCA <sub>AA</sub> GGUCCU5'

5'AAGAUGU<sup>A<span style="color:red">A</span>A</sup>CCAGGA3'
3'CUGCA <sub>AA</sub> GGUCCU5'

5'AAGAUGU<sup>AAA</sup>CCAGGA3'
3'CUGC<span style="color:red">A</span> <sub>AA</sub> GGUCCU5'

5'AAG<span style="color:red">A</span>UGU<sup>AAA</sup>CCAGGA
3'CUGCA <sub>AA</sub> GGUCCU

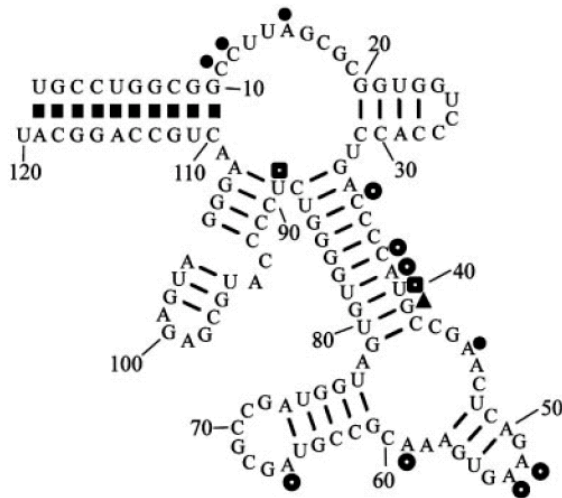5'AAGAUGU<sup>AAA</sup>CCAGG<span style="color:red">A</span>
3'CUGCA <sub>AA</sub> GGUCCU

5'AAGAUGU<sup>AAA</sup>CC<span style="color:red">A</span>GGA
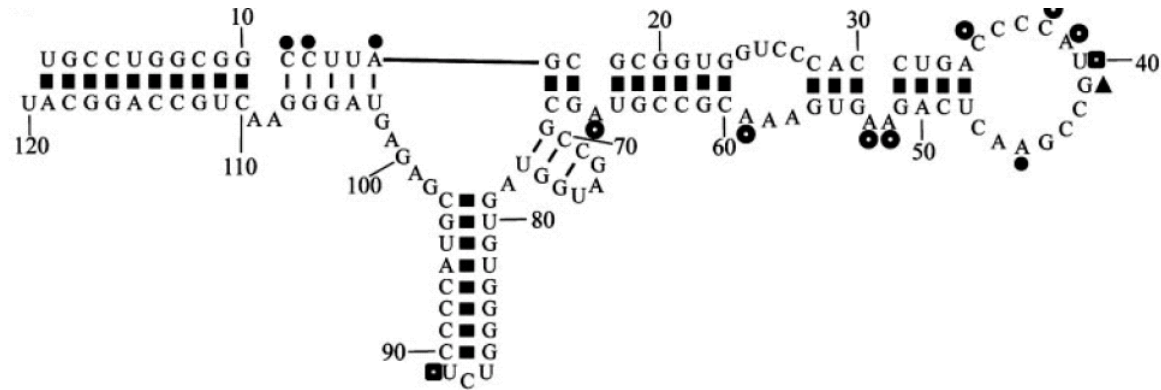3'CUGCA <sub>AA</sub> GGUCCU

# Including Results from Chemical Probing Improves Secondary Structure Prediction
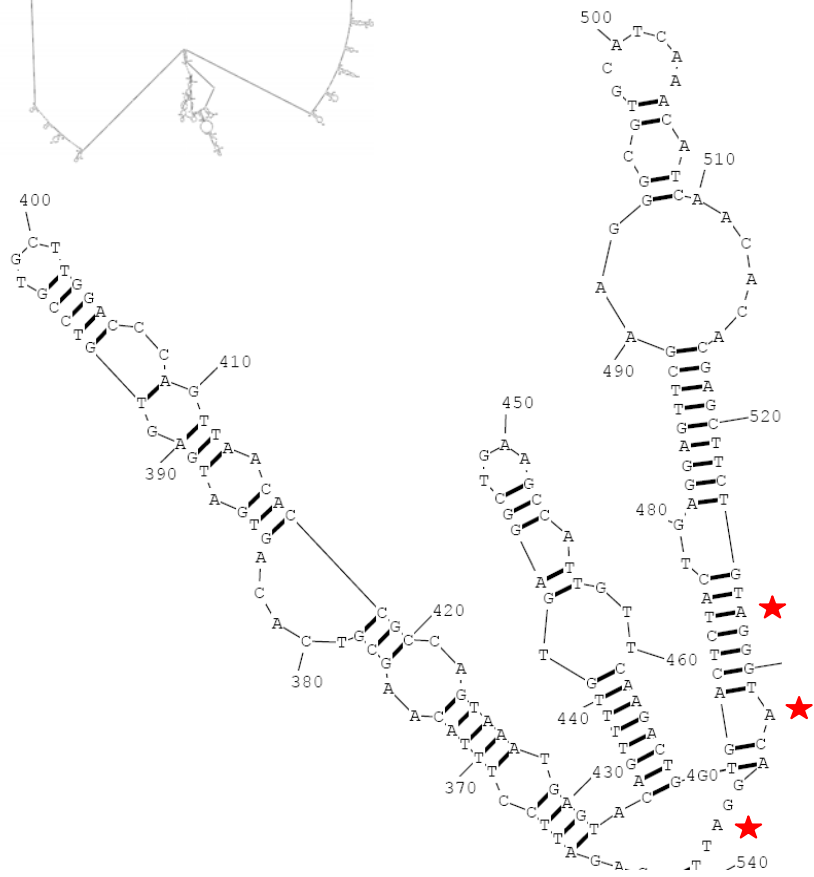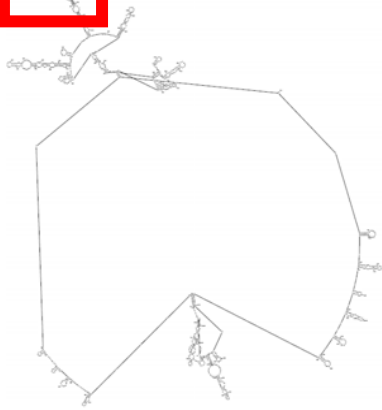
## *E.coli* 5S rRNA
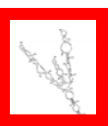


**LOWEST 26.3 %**
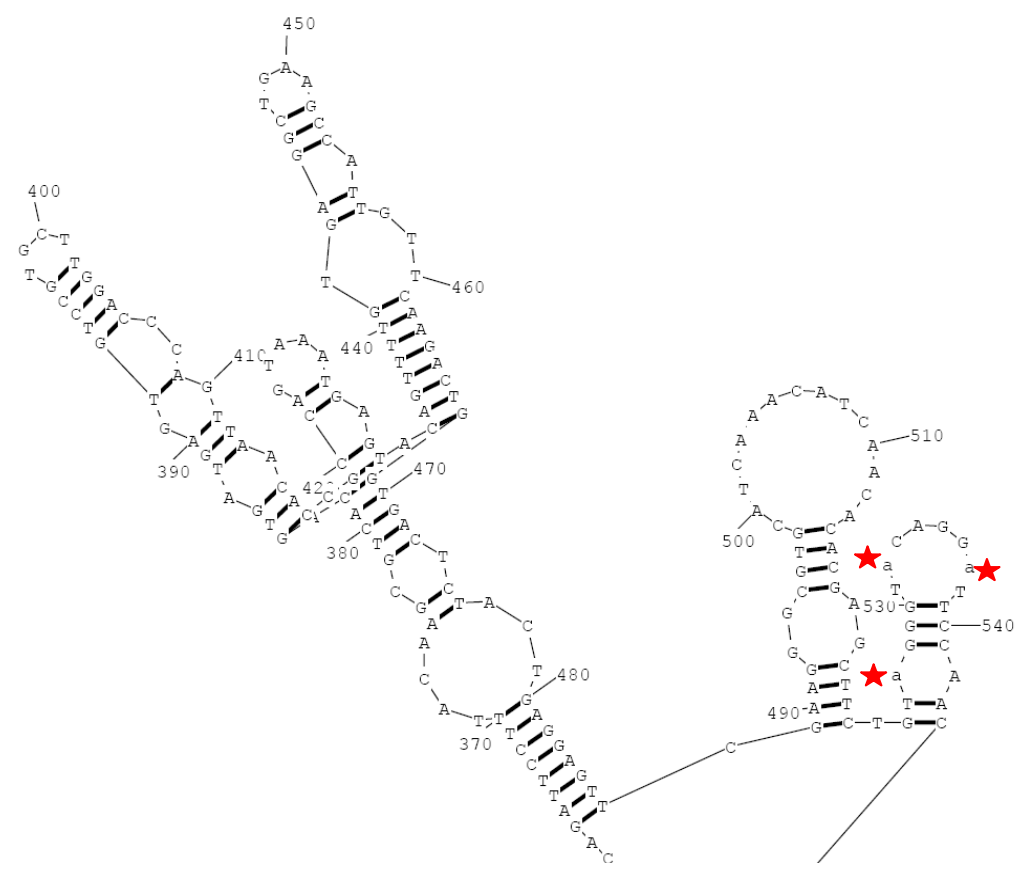**BEST        86.8 %**

**Folded with constraints**
**from *in vivo* chemical modification**
**LOWEST 86.8 %**
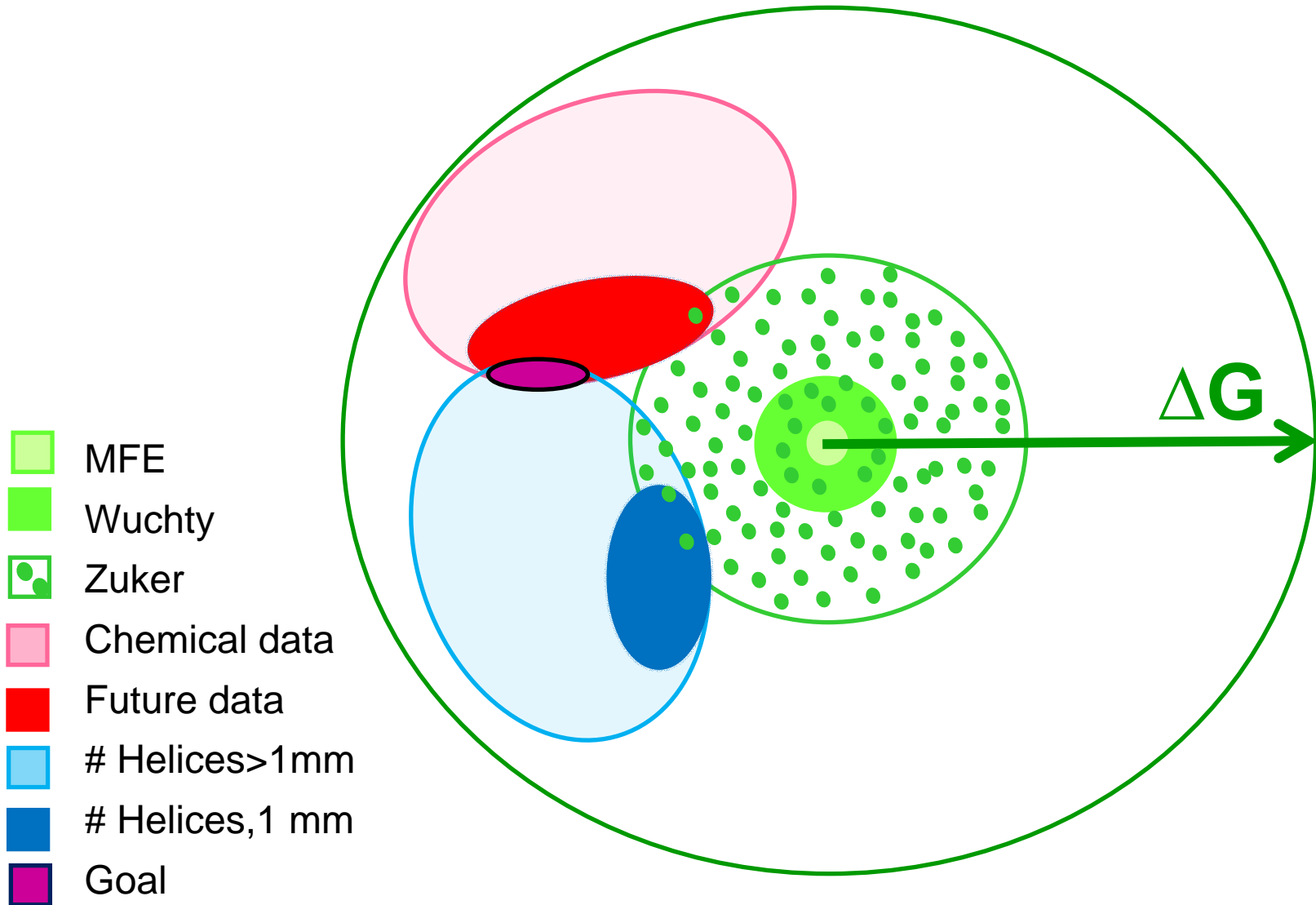**BEST        97.4 %**

# 3 Restraints Can Change the Lowest Energy Fold



Native -341.1 kcal/mol

A527, A532, A537
restrained to be single stranded
-341.0.kcal/mol

# Free Energy Landscape of STMV RNA

MFE

Wuchty

Zuker

Chemical data

Future data

# Helices>1mm

# Helices,1 mm

Goal

$\Delta G$

# How can Parallel Computing Help Solve the RNA Folding Problem?

- **Utilize tree structure of RNA secondary structure prediction**

- **Expand range of free energies that can be computed for an RNA free energy landscape**

- **Explore more possible RNA structures**

# Acknowledgements

- **Deb Mathews, University California Riverside**
- **David Mathews, University of Rochester**
- **Jeanmarie Verchot-Lubicz, Oklahoma State University**
- **Cal Lemke, Oklahoma University greenhouses**

- **Lab Members:**
  **Xiaobo Gu, Steven Harris, Koree Clanton-Arrowood, Brina Gendhar, Shelly Sedberry, Brian Doherty, Ted Gibbons, Sean Lavelle, John McGurk, Becky Myers, Mai-Thao Nguyen, Samantha Seaton, Jon Stone**
- **Funding**