

**CSE 527**

**Lecture 17, 11/24/04**

RNA Secondary Structure Prediction

# Outline

- What is it
- How is it Represented
- Why is it important
- Examples
- Approaches

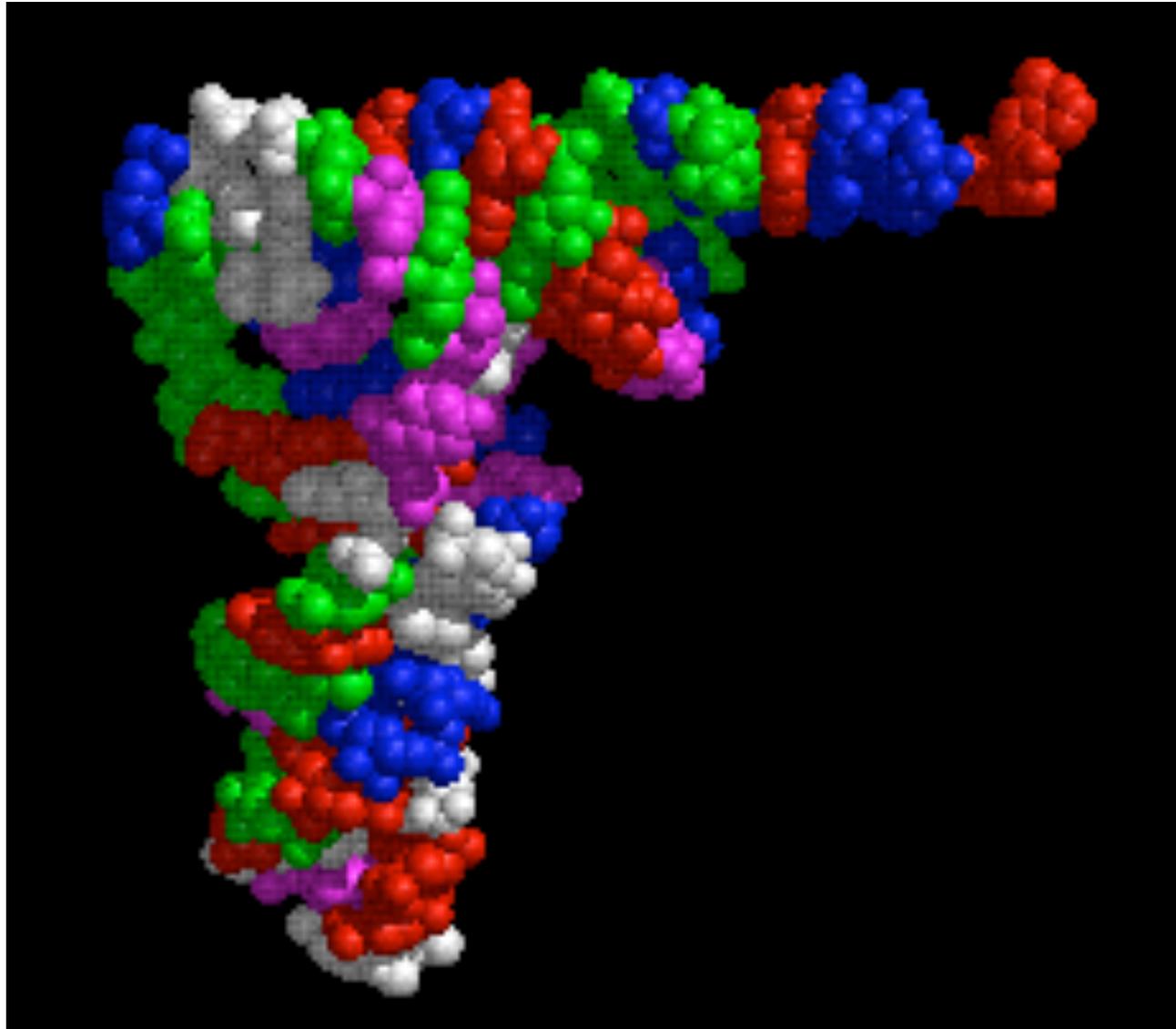
# RNA Structure

- Primary Structure: Sequence
- Secondary Structure: Pairing
- Tertiary Structure: 3D shape

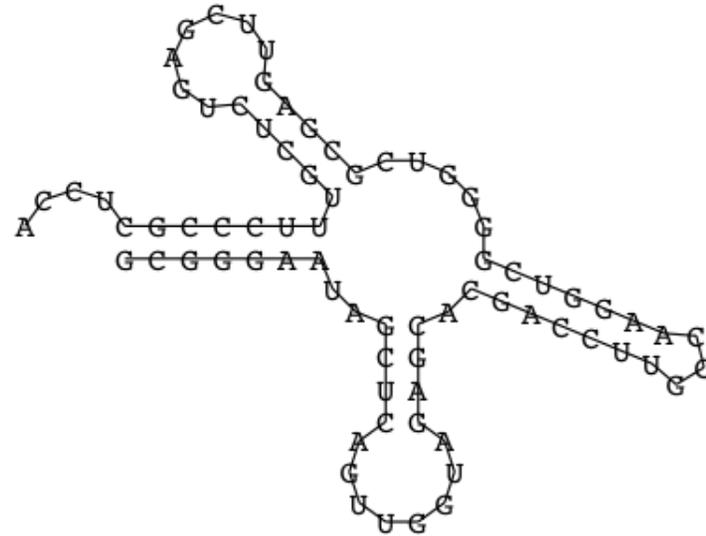
# RNA Pairing

- Watson-Crick Pairing
  - C - G ~ 3 kcal/mole
  - A - U ~ 2 kcal/mole
- “Wobble Pair” G - U ~ 1 kcal/mole
- Non-canonical Pairs (esp. if modified)

# A tRNA 3d Structure



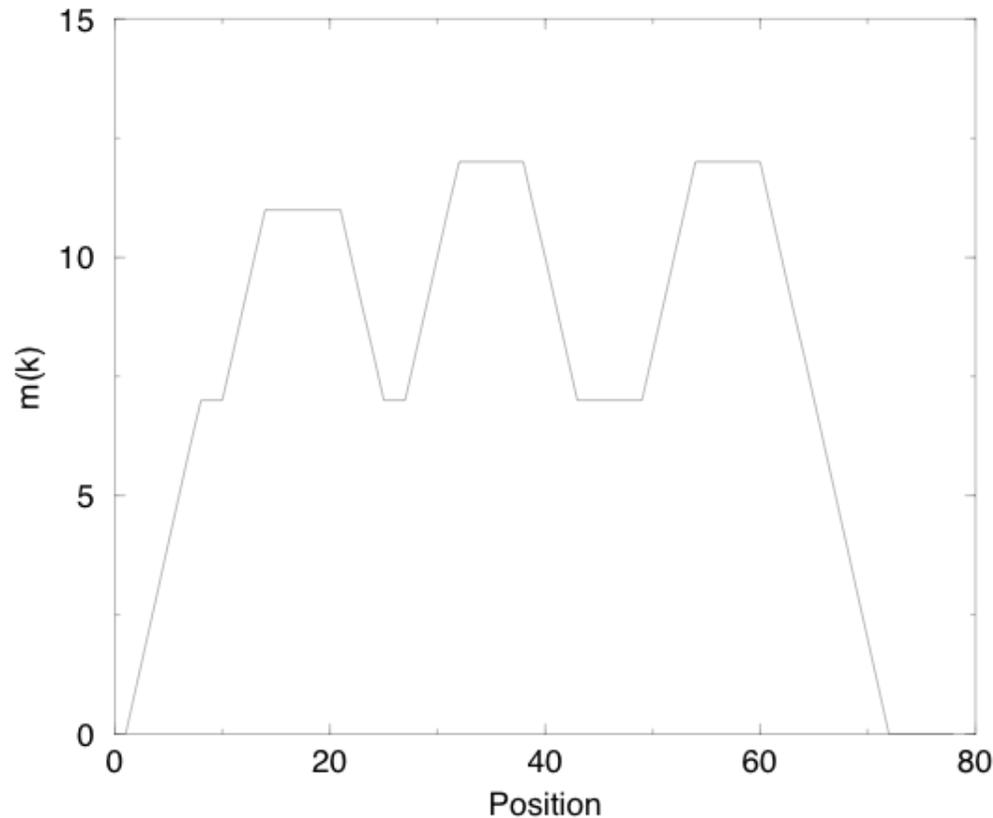
# tRNA - Alt. Representations



**Figure 1:** a) The spatial structure of the phenylalanine tRNA form yeast

b) The secondary structure extracts the most important information about the structure, namely the pattern of base pairings.

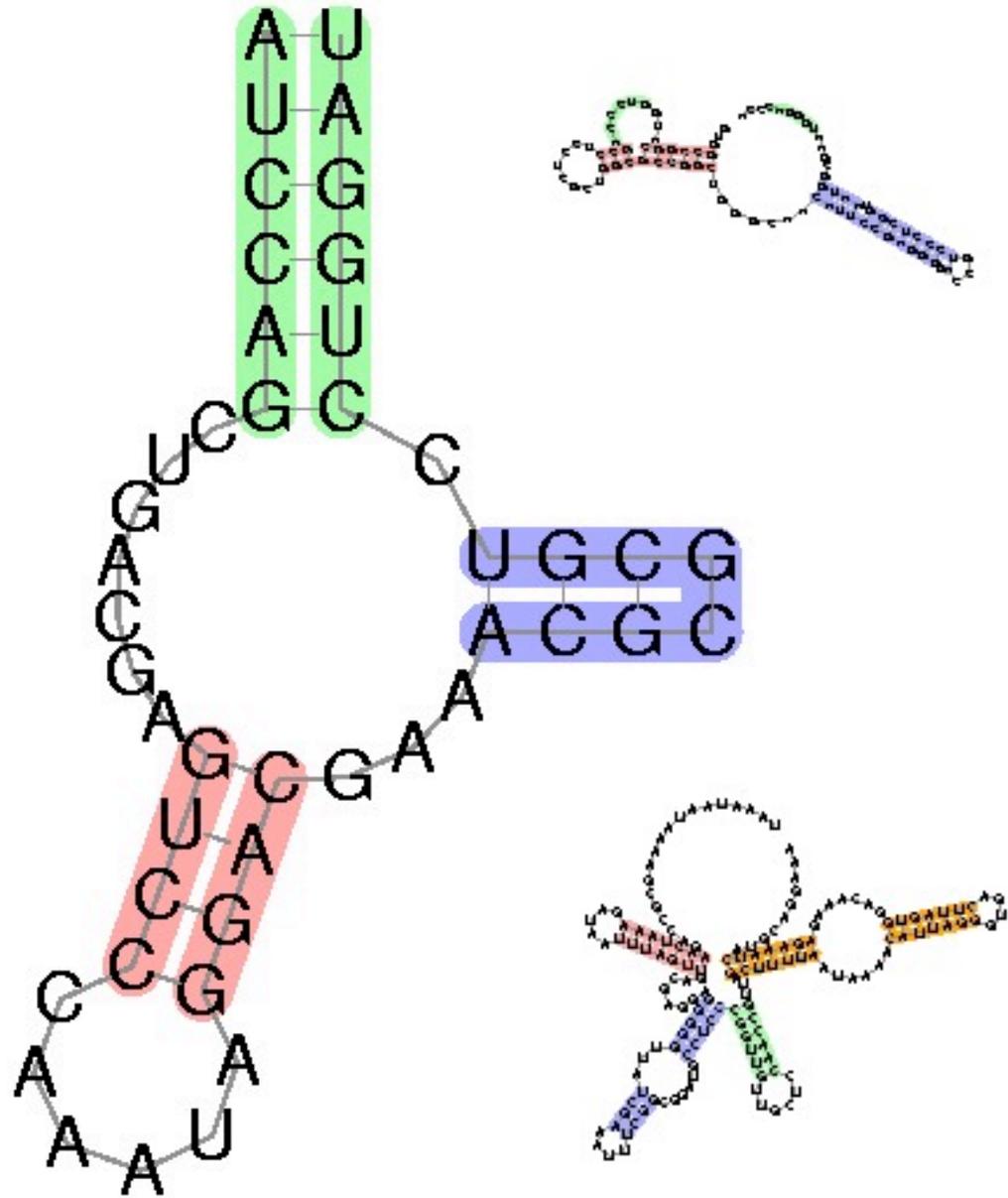
# A “Mountain” diagram



**Figure 3:** Mountain representation of the tRNA secondary structure shown in Figure 1. The three plateaus correspond to the three hairpin loops of the clover leaf structure.

# Why?

- RNA's fold, and function
- Nature uses what works



# Importance

- Ribozymes (RNA Enzymes)
- Retroviruses
- Effects on transcription, translation, splicing...
- Functional RNAs: rRNA, tRNA, snRNA, snoRNA, micro RNA, RNAi, riboswitches, regulatory elements in 3' & 5' UTRs, ...



# Definitions

- Sequence  $5' r_1 r_2 r_3 \dots r_n 3'$  in  $\{A, C, G, T\}$
- A **Secondary Structure** is a set of pairs  $i \bullet j$  s.t.

1.  $i < j-4$

2. if  $i \bullet j$  &  $i' \bullet j'$  are two pairs with  $i \leq i'$ , then

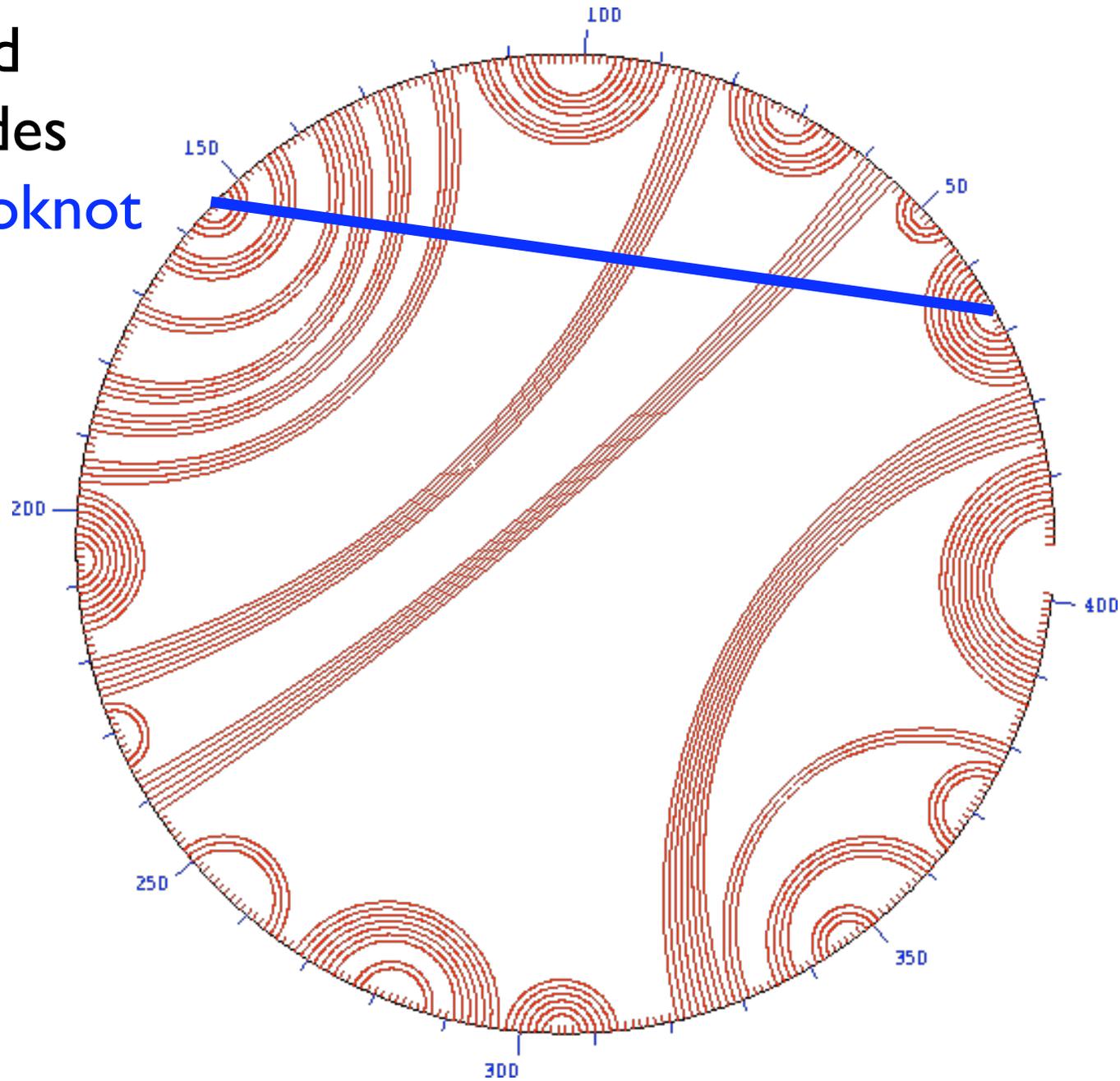
A.  $i = i' \text{ \& } j = j'$ , or

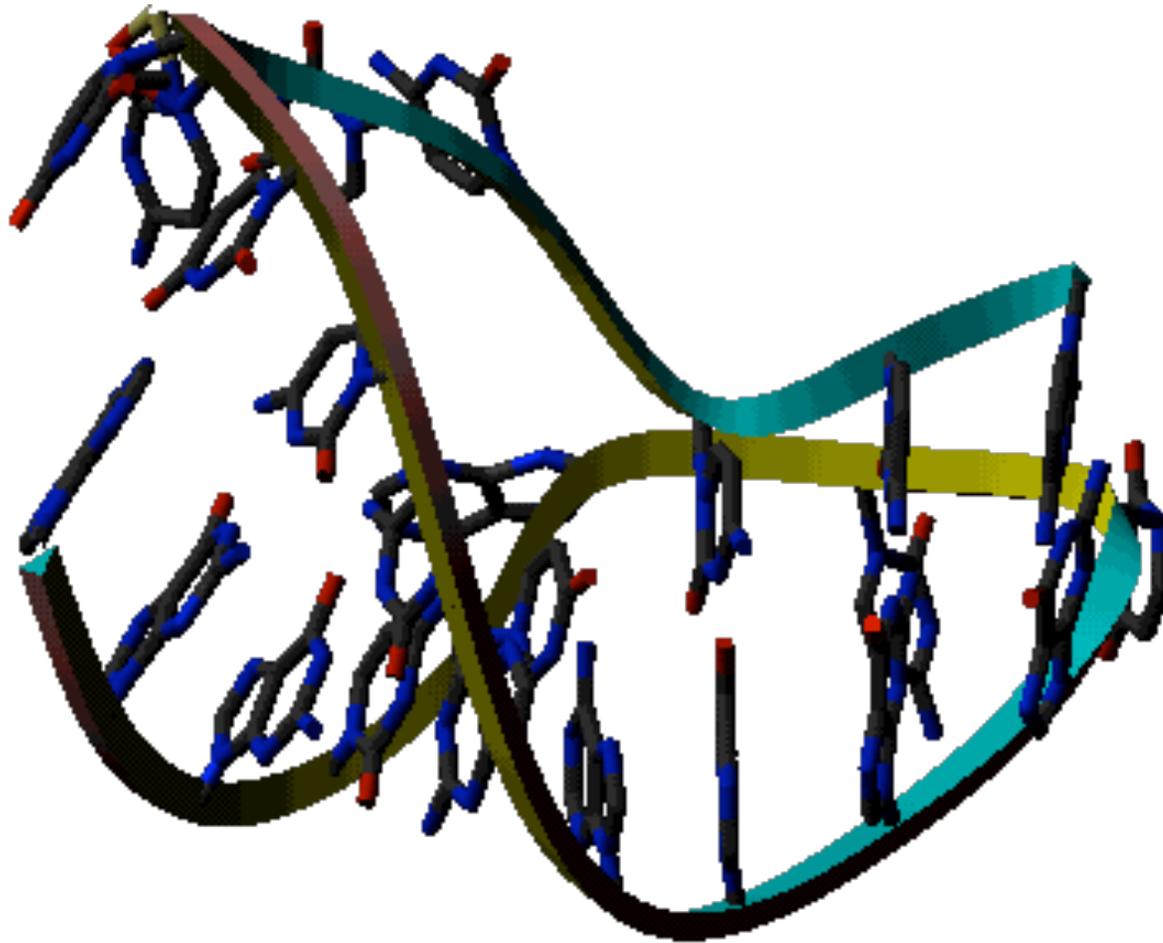
B.  $j < i'$ , or

C.  $i < i' < j' < j$

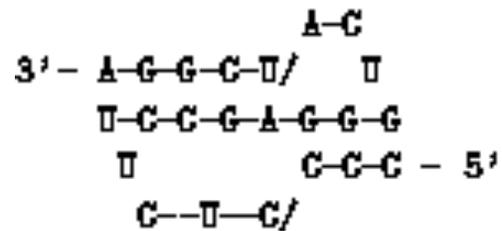
} First pair precedes 2nd,  
or is nested within it.  
No “pseudoknots.”

Nested  
Precedes  
Pseudoknot





The corresponding secondary structure is:



# A Pseudoknot

# Approaches to Structure Prediction

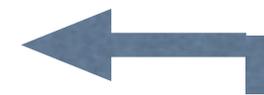
- Maximum Pairing
  - + simple
  - too inaccurate
- Minimum Energy
  - + Works on single sequences
  - Ignores pseudoknots
  - Only finds “optimal” fold
- Partition Function
  - + Finds all folds
  - Ignores pseudoknots

# Approaches, II

- Comparative sequence analysis
  - + handles all pairings (incl. pseudoknots)
  - requires several (many?) aligned, appropriately diverged
- Stochastic Context-free Grammars
  - Roughly combines min energy & comparative, but no pseudoknots
- Physical experiments (x-ray crystallography, NMR)

# Nussinov: Max Pairing

- $B(i,j)$  = # pairs in optimal pairing of  $r_i \dots r_j$
- $B(i,j) = 0$  for all  $i, j$  with  $i \geq j-4$ ; otherwise
- $B(i,j) = \max$  of:
  1.  $B(i+1,j)$
  2.  $B(i,j-1)$
  3.  $B(i+1,j-1)$  +(if  $r_i$  pairs with  $r_j$  then 1 else 0)
  4.  $\max \{ B(i,k)+B(k+1,j) \mid i < k < j \}$



Time:  $O(n^3)$

# Pair-based Energy Minimization

- $E(i,j)$  = energy of pairs in optimal pairing of  $r_i \dots r_j$
- $E(i,j) = \infty$  for all  $i, j$  with  $i \geq j-4$ ; otherwise
- $E(i,j) = \min$  of:

1.  $E(i+1, j)$

2.  $E(i, j-1)$

3.  $E(i+1, j-1) + e(r_i, r_j)$

4.  $\min \{ E(i, k) + E(k+1, j) \mid i < k < j \}$



energy of one pair

Time:  $O(n^3)$



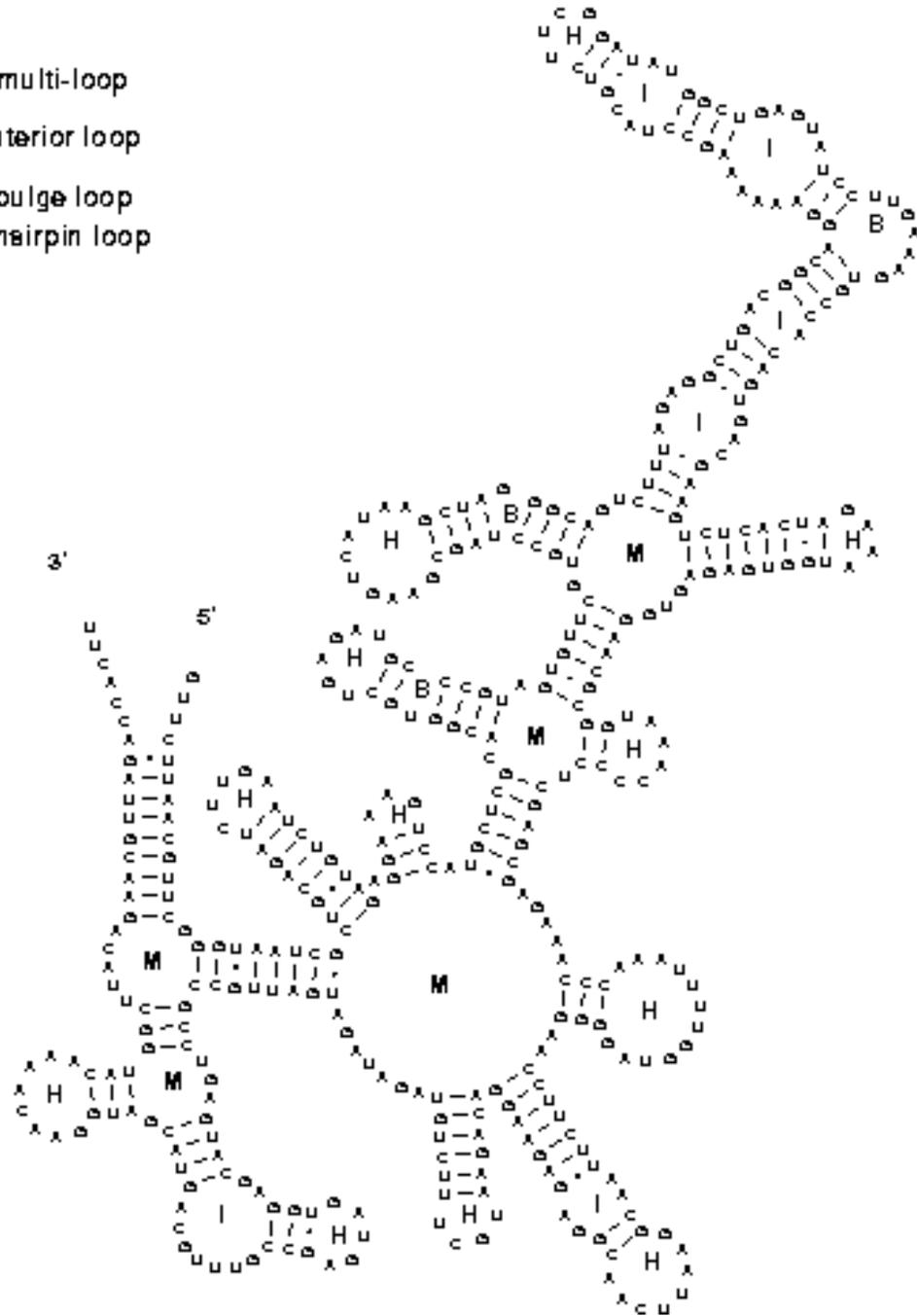
# Loop-based Energy Minimization

- Detailed experiments show that it's more accurate to model based on loops, rather than just pairs
- Loop types
  - Stack
  - Hairpin loop
  - Bulge
  - Interior loop

*Bacillus subtilis RNase P RNA*

- M** - multi-loop
- I** - interior loop
- B** - bulge loop
- H** - hairpin loop

# Loop Examples



# Zuker: Loop-based Energy, I

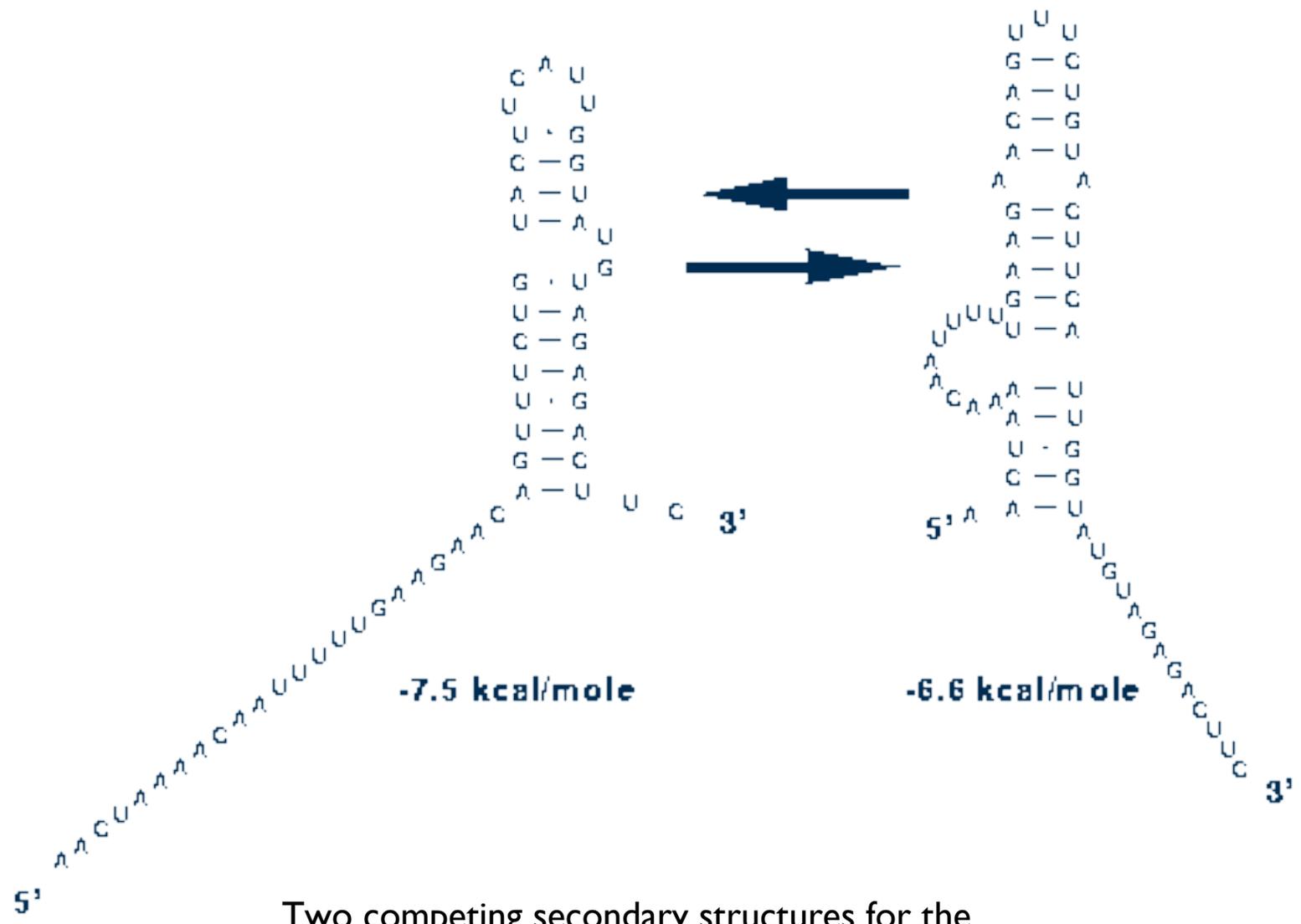
- $W(i,j)$  = energy of optimal pairing of  $r_i \dots r_j$
- $V(i,j)$  = as above, but forcing pair  $i \bullet j$
- $W(i,j) = V(i,j) = \infty$  for all  $i, j$  with  $i \geq j-4$
- $W(i,j) = \min(W(i+1,j), W(i,j-1), V(i+1,j-1), \min \{ E(i,k)+E(k+1,j) \mid i < k < j \} )$

# Zuker: Loop-based Energy, II

- $V(i,j) = \min(\overset{\text{hairpin}}{eh(i,j)}, \overset{\text{stack}}{es(i,j)} + V(i+1, j-1), \overset{\text{bulge/interior}}{VBI(i,j)}, \overset{\text{multi-loop}}{VM(i,j)})$
  - $VM(i,j) = \min \{ W(i,k) + W(k+1, j) \mid i < k < j \}$
  - $VBI(i,j) = \min \{ ebi(i,j,i',j') + V(i', j') \mid i < i' < j' < j \ \& \ i' - i + j - j' > 2 \}$
- Time:  $O(n^4)$   
 $O(n^3)$  possible if  $ebi(\cdot)$  is “nice”

# Suboptimal Energy

- There are always alternate folds with near-optimal energies. Thermodynamics predicts that populations of identical molecules will exist in different folds; individual molecules even flicker among different folds
- Zuker's algorithm can be modified to find suboptimal folds
- McCaskill gives a more elaborate dynamic programming algorithm calculating the "partition function," which defines the probability distribution over all these states.



Two competing secondary structures for the *Leptomonas collosoma* spliced leader mRNA.

# Example of suboptimal folding

Black dots:  
pairs in opt fold

Colored dots:  
pairs in folds  
2-5% worse  
than  
optimal fold

