

RNA

secondary structure prediction and analysis

Resources

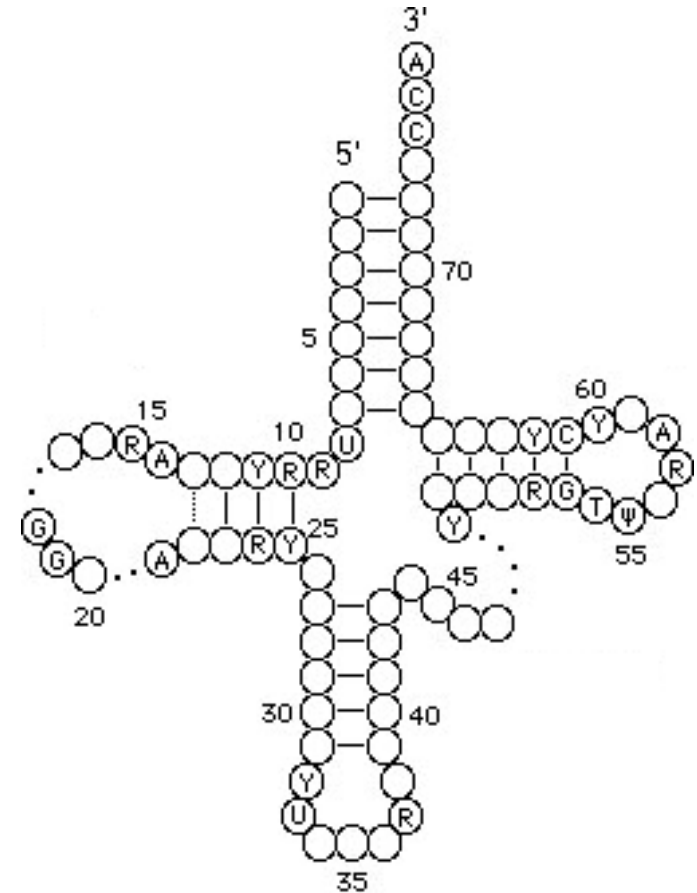
- Lecture Notes from previous years: Takis Benos
- Covariance algorithm: Eddy and Durbin, Nucleic Acids Research, v22: 11, 2079
- Useful lecture slides from Larry Ruzzo, U of Washington and Phillip Compeau, UCSD

Outline

- RNA Folding
- Dynamic Programming for RNA secondary structure prediction
 - Nussinov et al and Zuker et al algorithms
- Covariance Model
 - Eddy and Durbin

Various types of RNA

- messenger RNA (mRNA)
- transfer RNA (tRNA)
- Ribosomal RNA (rRNA)
- small interfering RNA (siRNA)
- micro RNA (miRNA)
- small nuclear RNA (snRNA)
- small nucleolar RNA (snoRNA)
- guide RNA (gRNA)
- efference RNA(eRNA)



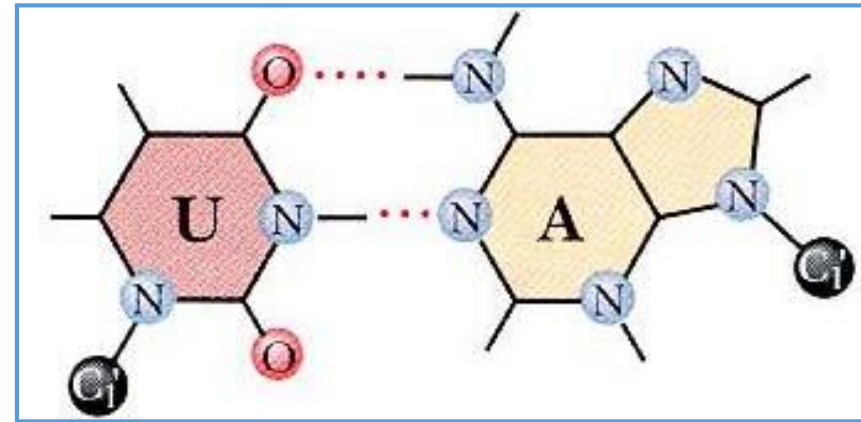
Non coding RNA (ncRNA)

- RNA that isn't translated into protein
- Includes:
 - tRNA, rRNA, snRNA, snoRNA, miRNA, gRNA, eRNA, pRNA, tmRNA
- What about mRNA?
 - 5' methylated cap
 - 5'-UTR (un-translated regions)
 - CDS (coding sequence)
 - 3'-UTR
 - Poly-A tail
- mRNA contains untranslated regions (5'UTR, 3'UTR), but UTRs are not considered ncRNA

RNA Basics

- RNA bases: A, C, G, U
- Watson-Crick Pair
 - A-U (~ 2 kcal/mol)
 - G-C (~ 3 kcal/mol)
- Wobble pair
 - G-U (~ 1 kcal/mol)
- Non-Canonical pairs (modified suitably)
- Bases can only pair with one other base

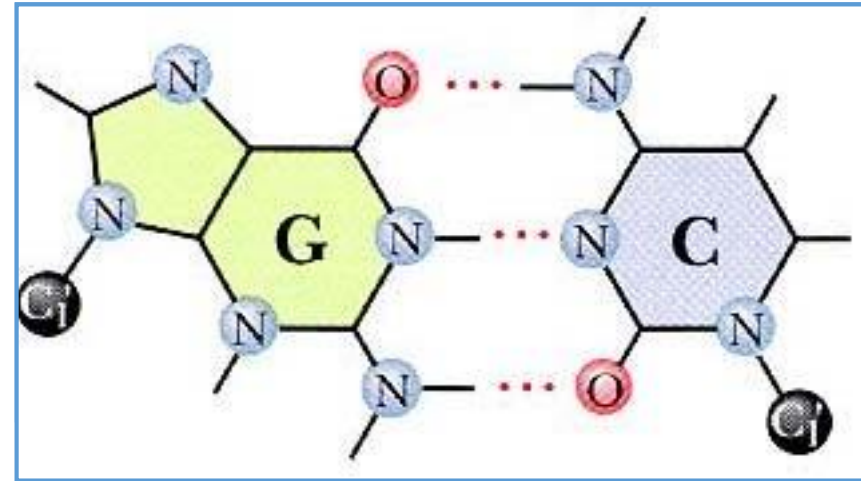
Two hydrogen bonds



RNA Basics

- RNA bases: A, C, G, U
- Canonical Base Pairs
 - A-U
 - G-C
 - G-U
- Bases can only pair with one other base

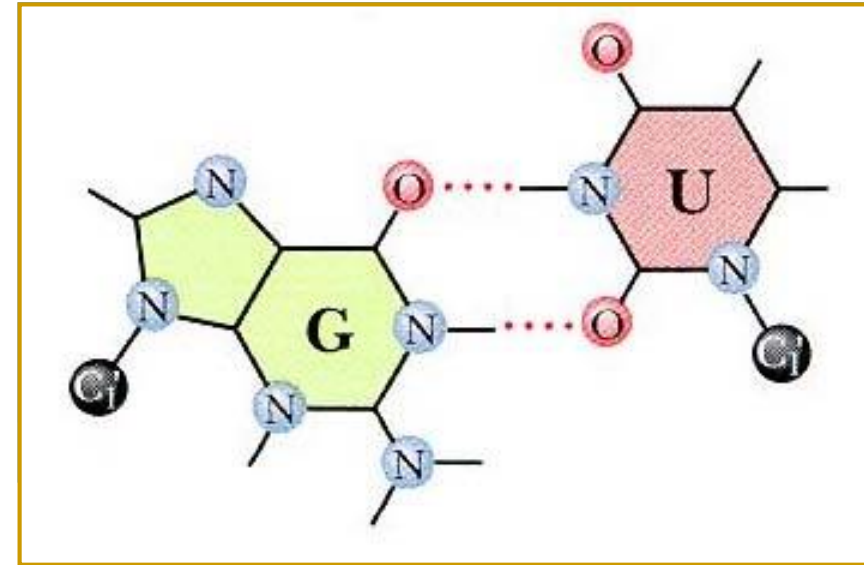
Three hydrogen bonds => more stable



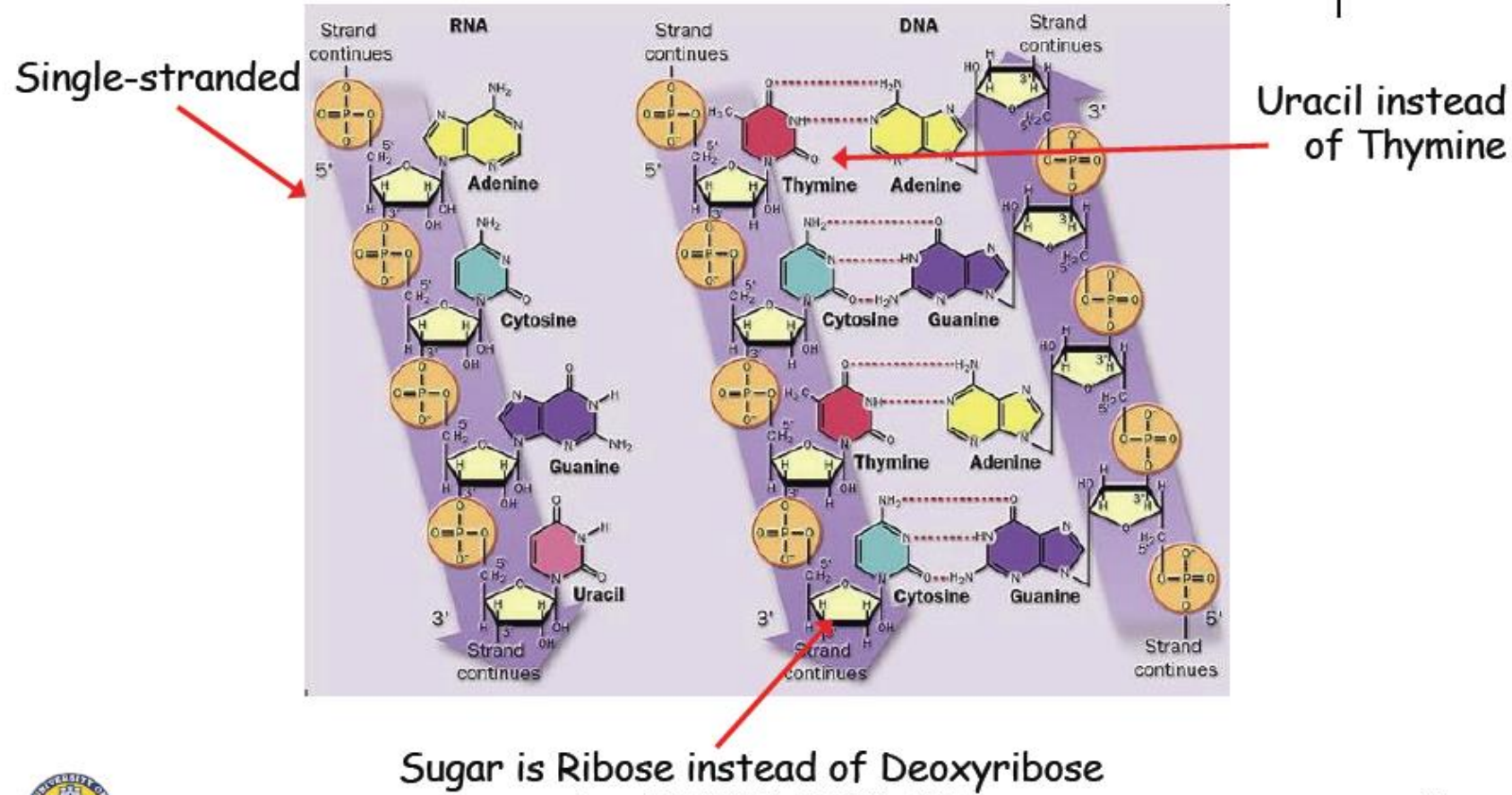
RNA Basics

- RNA bases: A, C, G, U
- Canonical Base Pairs
 - A-U
 - G-C
 - G-U
- Bases can only pair with one other base

‘wobble pairing’



How is RNA different from DNA?



Benos MSCBIO2070 25-27.Mar.2012

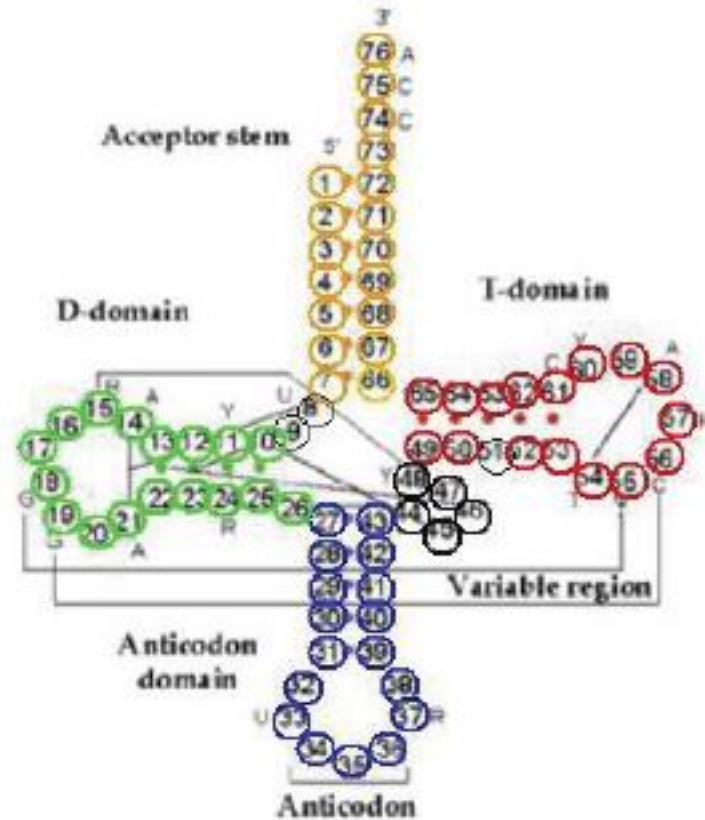
15

Secondary and Tertiary Structure

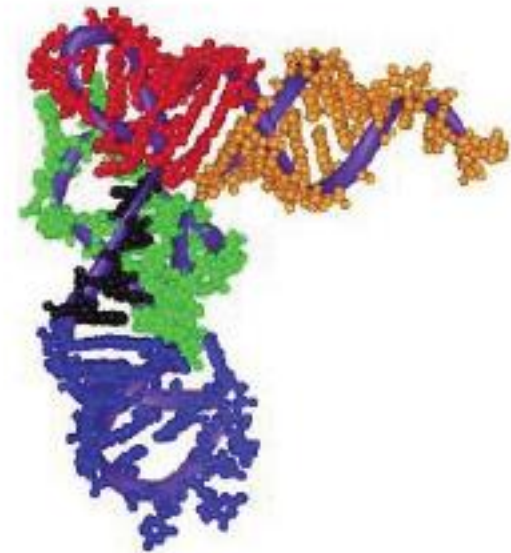


GCGGAUUUAGCUCAGUUGG
GAGAGCGCCAGACUGAAGA
UCUGGAGGUCCUGUGUUCG
AUCCACAGAAUUCGCACCA

**Primary
Structure**

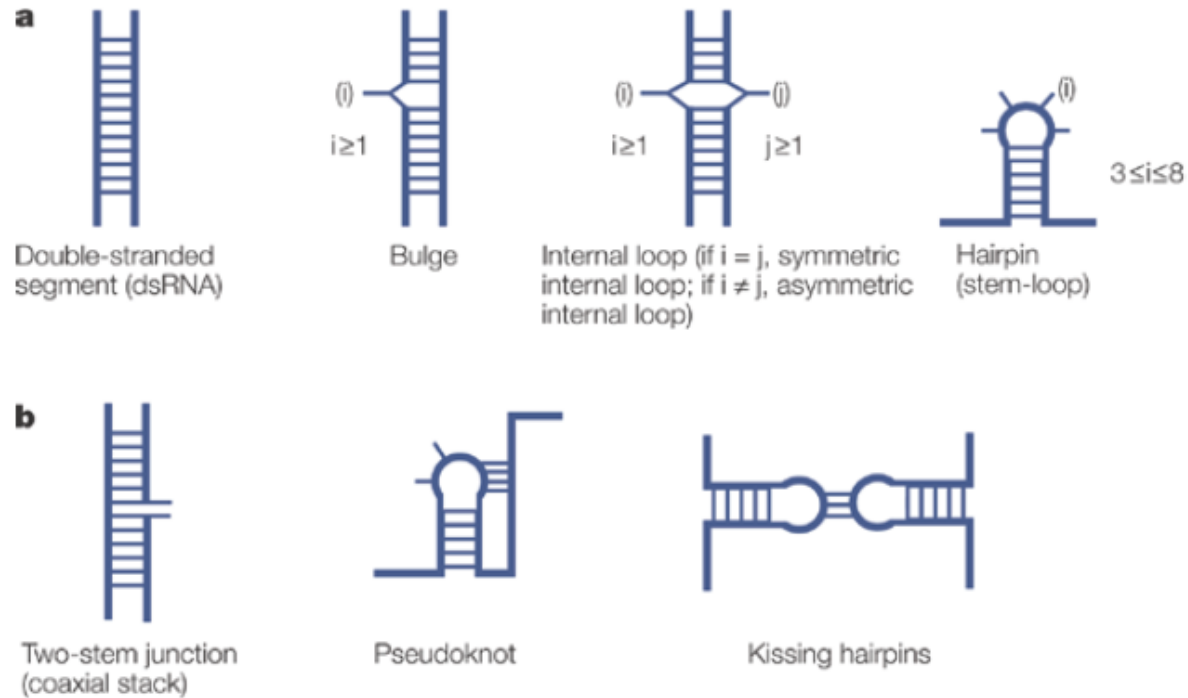


**Secondary
Structure**

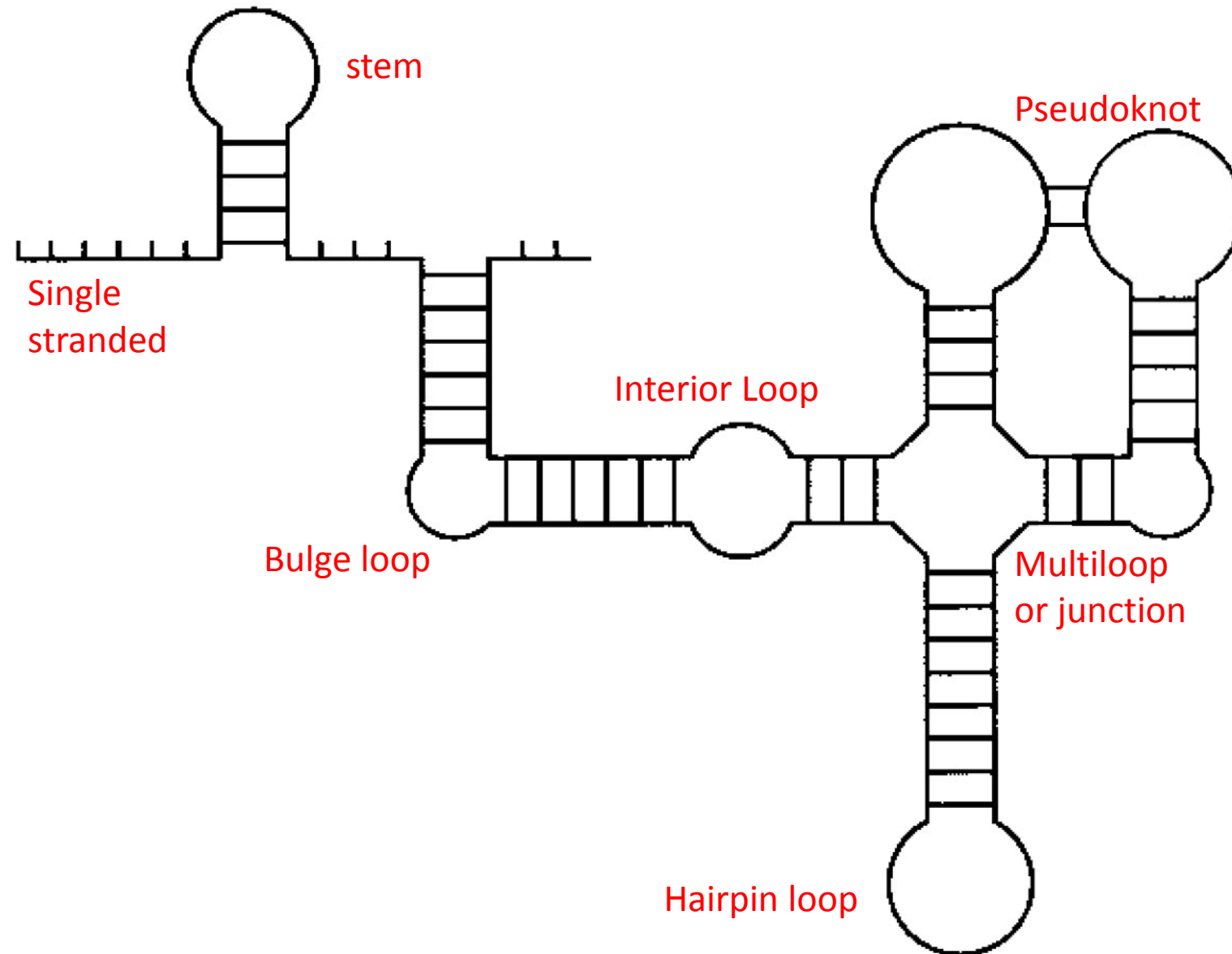


**Tertiary
Structure**

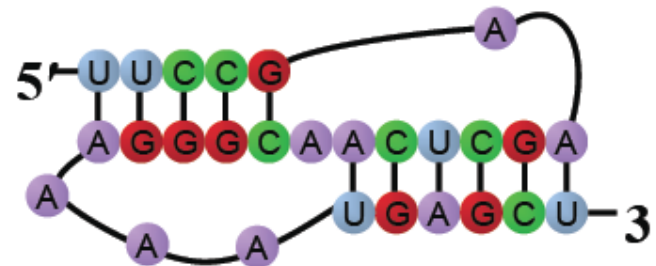
Secondary structural elements



RNA Secondary Structure/Motifs



Pseudoknots



bases pairs between a loop and positions outside the enclosing stem

Challenging to deal with them; however, the total number of pseudoknotted base pairs is relatively small

i.e. in *E. coli* SSU rRNA, 447 base pairs, only 8 are in pseudoknot structures



tRNA - Alt. Representations

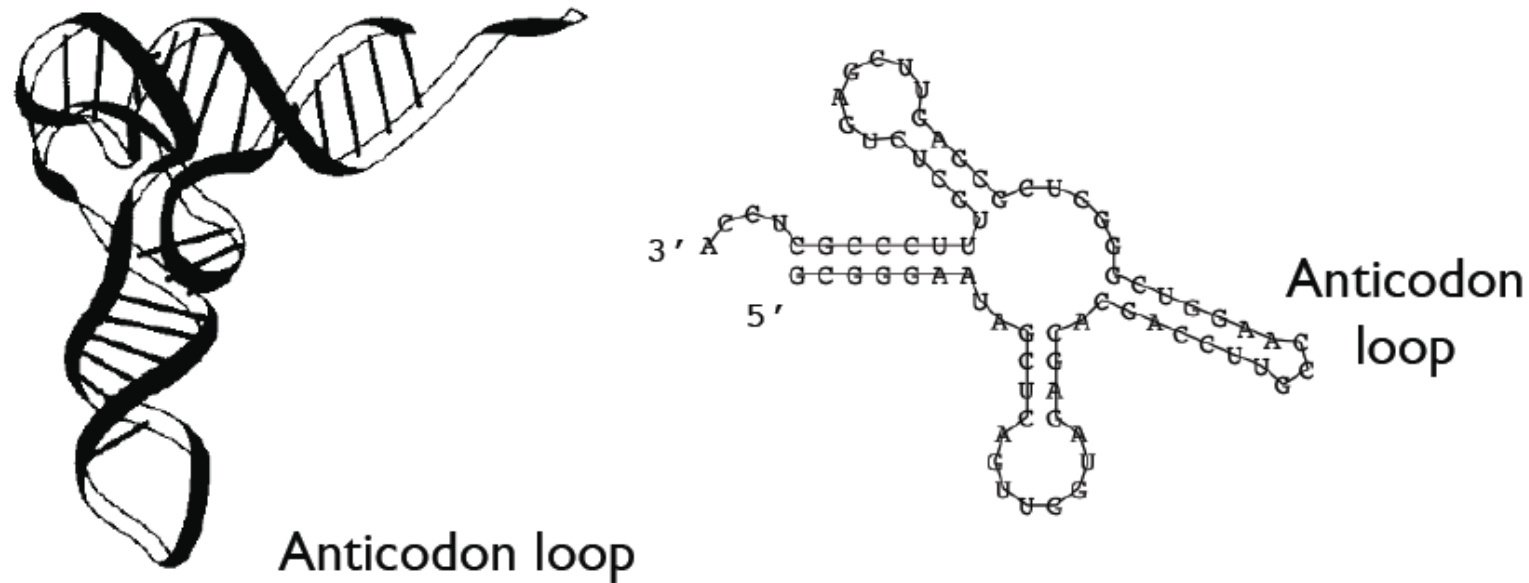
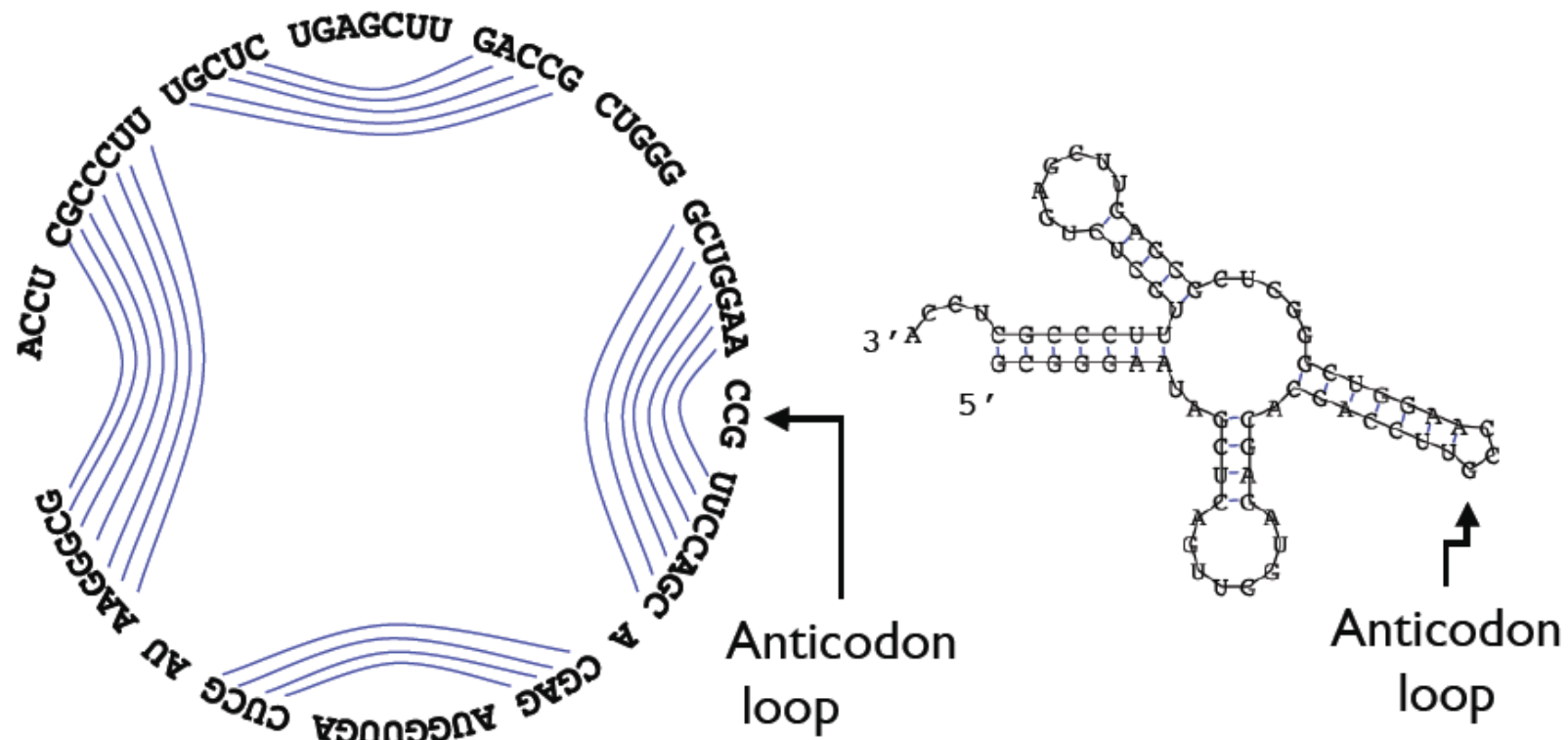


Figure 1: a) The spatial structure of the phenylalanine tRNA from yeast

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

tRNA - Alt.

Representations



RNA importance

- Ribozymes (RNA enzymes)
- Retroviruses
- Effects on transcription, translation, splicing..
- Functional RNAs: rRNA, tRNA, snRNA, snoRNA, microRNA, RNAi, riboswitches, regulatory elements in 3' and 5' UTR

RNA motifs

- Function depends on sequence and secondary structure
- Functions include
 - Protein binding
 - Basepairing to another RNA
 - Modifying a nucleic acid bond
- Types:
 - Single-strand regions
 - Helices (or stems)
 - Bulges
 - Hairpin loops
 - Internal loops
 - Junctions

RNA Motifs Regulatory Effects

- Regulations of translations
- Processing of RNA
- Catalytic modification of other RNAs
- Transport and position in the cell
- Stability of RNA-transcript
- Expression of encoded proteins

Why predict structures?

- Knowing the shape of a biomolecule is invaluable in drug design and understanding disease mechanisms
- Current physical methods (X-Ray, NMR) are too expensive and time-consuming
- Predict shape from sequence of bases
- Four basic structures: helices, loops, bulges and junctions

RNA secondary structure

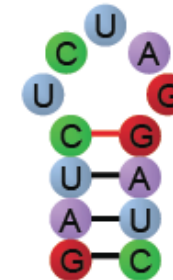
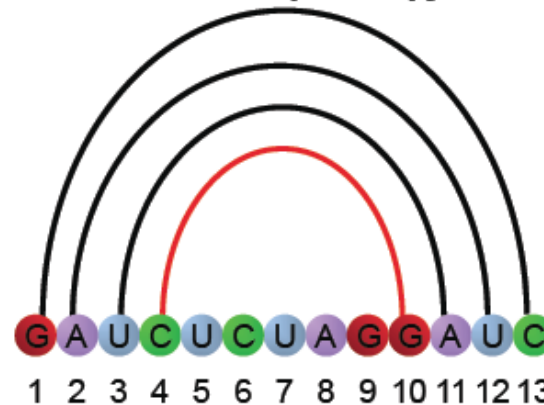
- What makes RNA fold?
- Problem: given an RNA sequence, find the set of base pairs that is “correct” or “optimal”
 - Maximize number of base pairs (Nussinov et al)
 - Minimize energy (Zucker et al)
- Search problem: very high number of possible structures
- Algorithm: dynamic programming
 - Cannot handle pseudoknots



Structure Representation

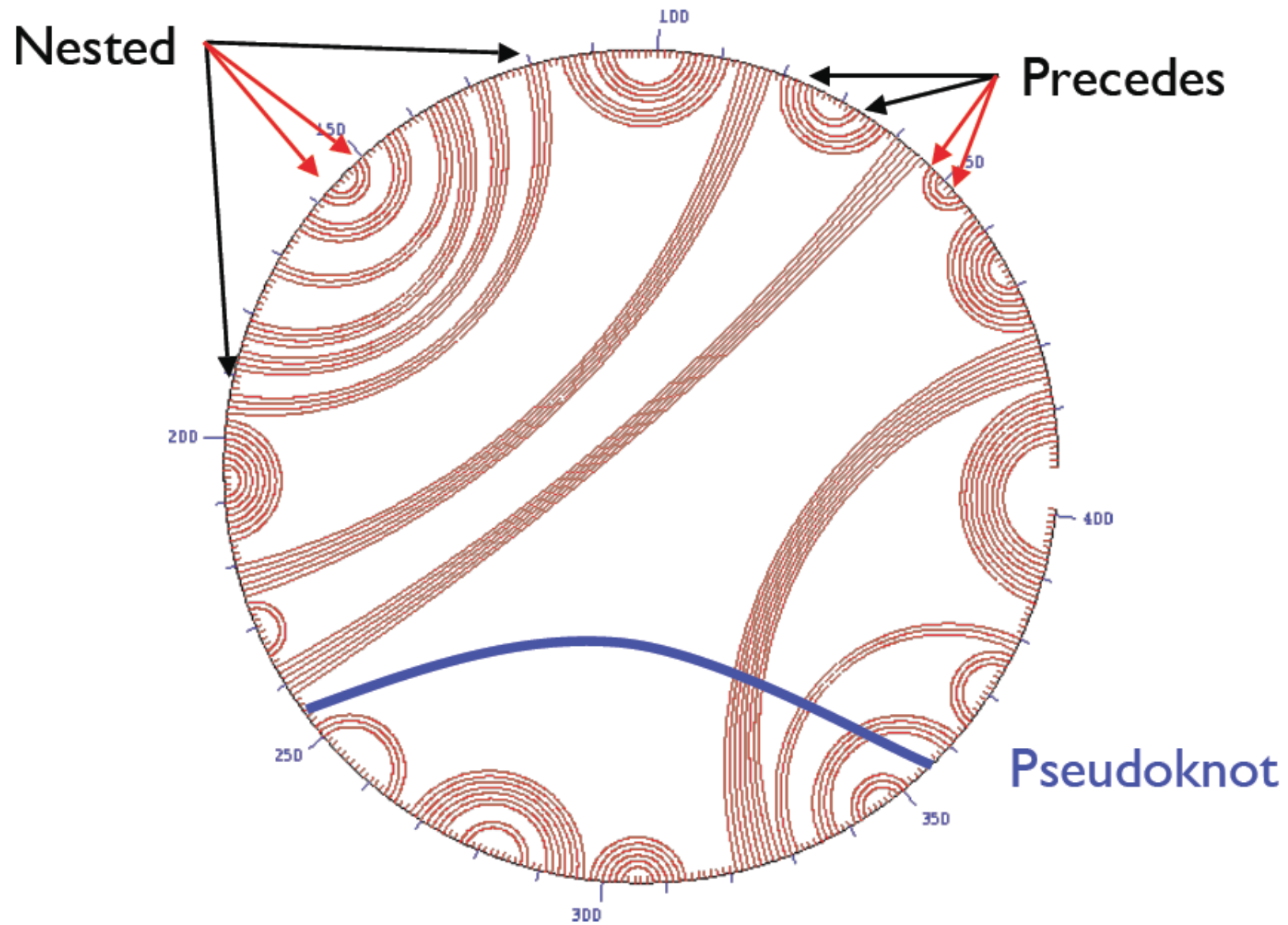
- secondary structure described as a graph
- base pairs are described via pairs of indices (i, j), indicating links between base vertices

$S = \{(1,13), (2,12), (3,11), (4,10)\}$

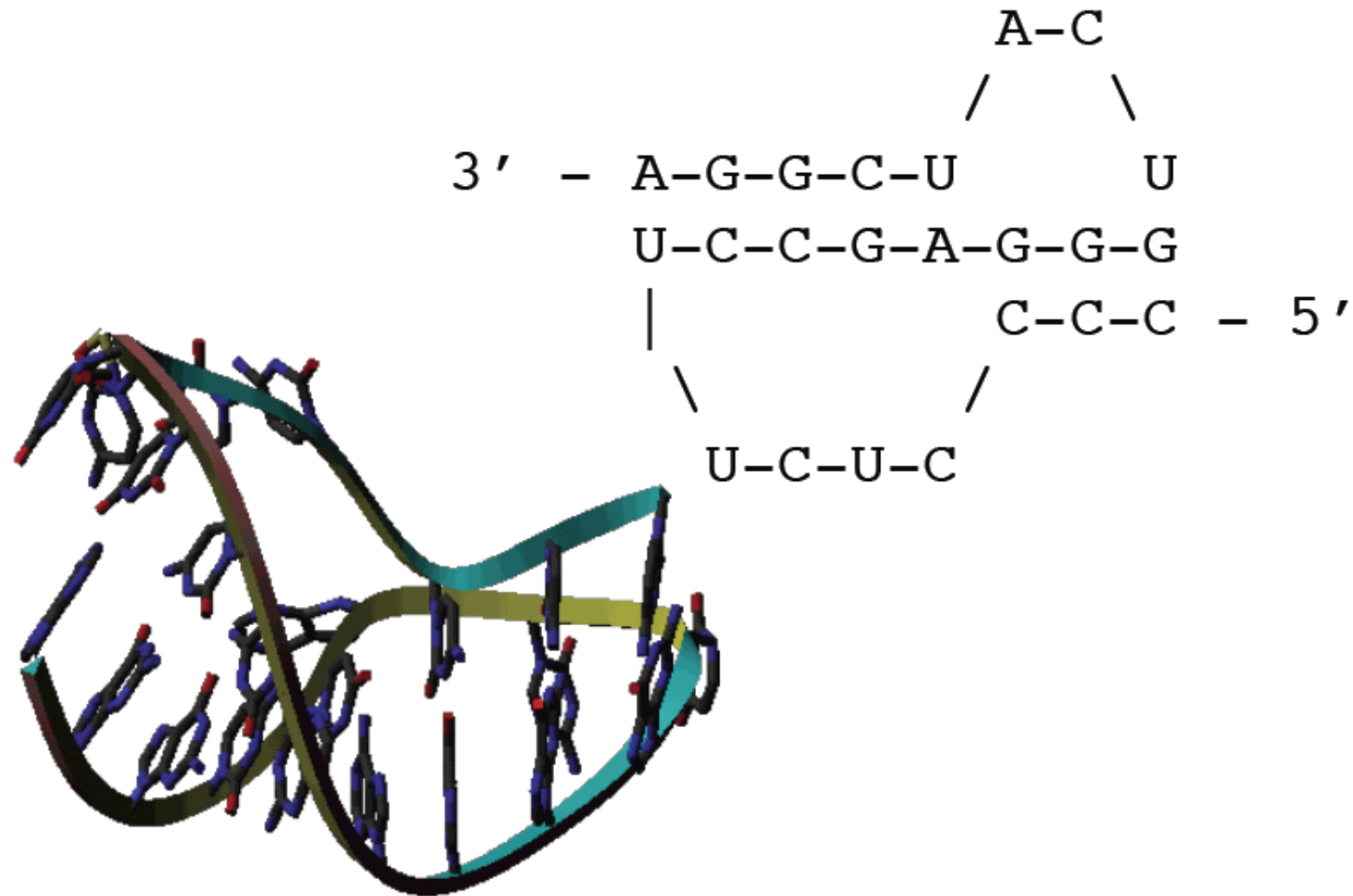


Definitions

- 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. $i = i'$ & $j = j'$, or
 - B. $j < i'$, or
 - C. $i < i' < j' < j$
- } First pair precedes 2nd, or is nested within it. No “pseudoknots.”



A Pseudoknot



Approaches to Structure Prediction

- Maximum Pairing
 - + works on single sequences
 - + 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 sequences
- Stochastic Context-free Grammars
 - Roughly combines min energy & comparative, but no pseudoknots
- Physical experiments (x-ray crystallography, NMR)

Base Pair Maximization: Dynamic Programming

- $S(i, j)$ is the folding of the RNA subsequence of the strand from index i to index j which results in the highest number of base pairs.
- Recurrence:

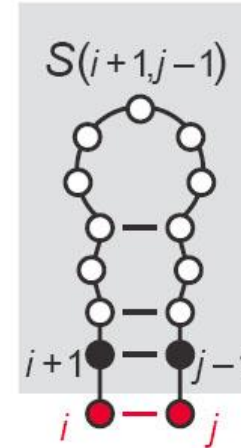
Base Pair Maximization: Dynamic Programming

- $S(i, j)$ is the folding of the RNA subsequence of the strand from index i to index j which results in the highest number of base pairs.
- Recurrence:

$$S(i, j) = \max \left\{ \begin{array}{l} \\ \\ \\ \end{array} \right.$$

Base Pair Maximization: Dynamic Programming

- $S(i, j)$ is the folding of the RNA subsequence of the strand from index i to index j which results in the highest number of base pairs.



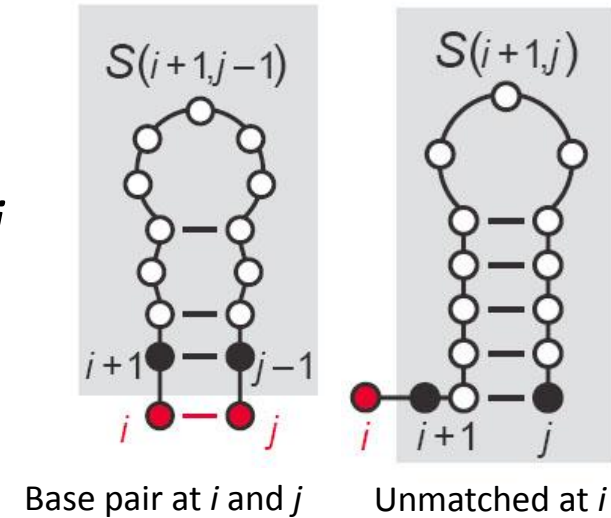
Base pair at i and j

- Recurrence:

$$S(i, j) = \max \left\{ \begin{array}{l} S(i+1, j-1) + 1 \quad (\text{if } (i, j) \text{ base pair}) \\ \vdots \end{array} \right.$$

Base Pair Maximization: Dynamic Programming

- $S(i, j)$ is the folding of the RNA subsequence of the strand from index i to index j which results in the highest number of base pairs.

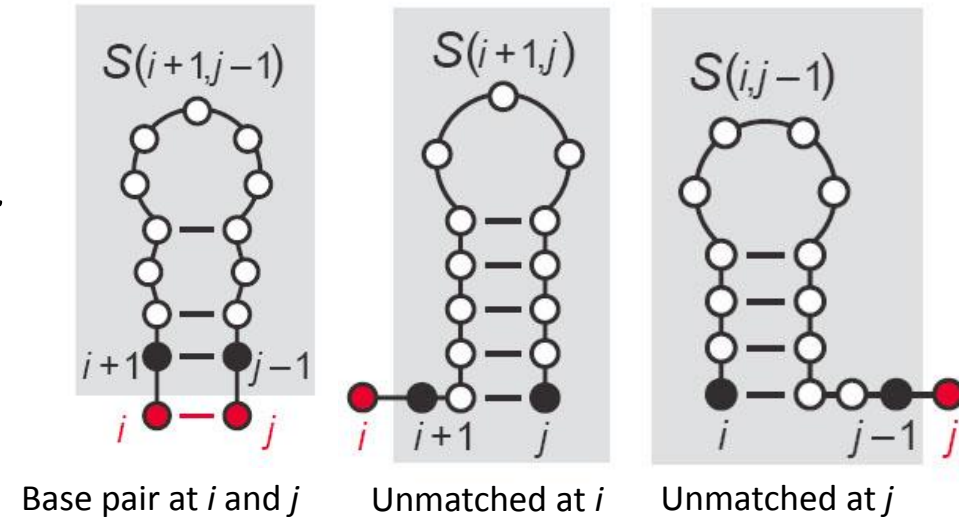


- Recurrence:

$$S(i, j) = \max \begin{cases} S(i+1, j-1) + 1 & \text{(if } (i, j) \text{ base pair)} \\ S(i+1, j) \\ \end{cases}$$

Base Pair Maximization: Dynamic Programming

- $S(i, j)$ is the folding of the RNA subsequence of the strand from index i to index j which results in the highest number of base pairs.

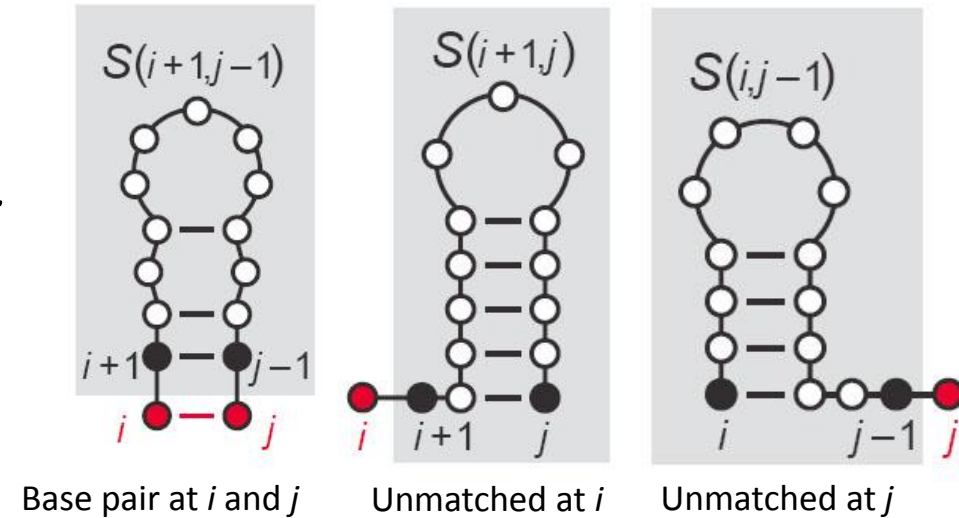


- Recurrence:

$$S(i, j) = \max \begin{cases} S(i+1, j-1) + 1 & \text{(if } (i, j) \text{ base pair)} \\ S(i+1, j) \\ S(i, j-1) \end{cases}$$

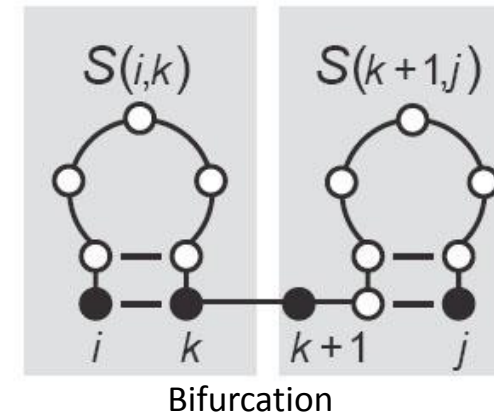
Base Pair Maximization: Dynamic Programming

- $S(i, j)$ is the folding of the RNA subsequence of the strand from index i to index j which results in the highest number of base pairs.



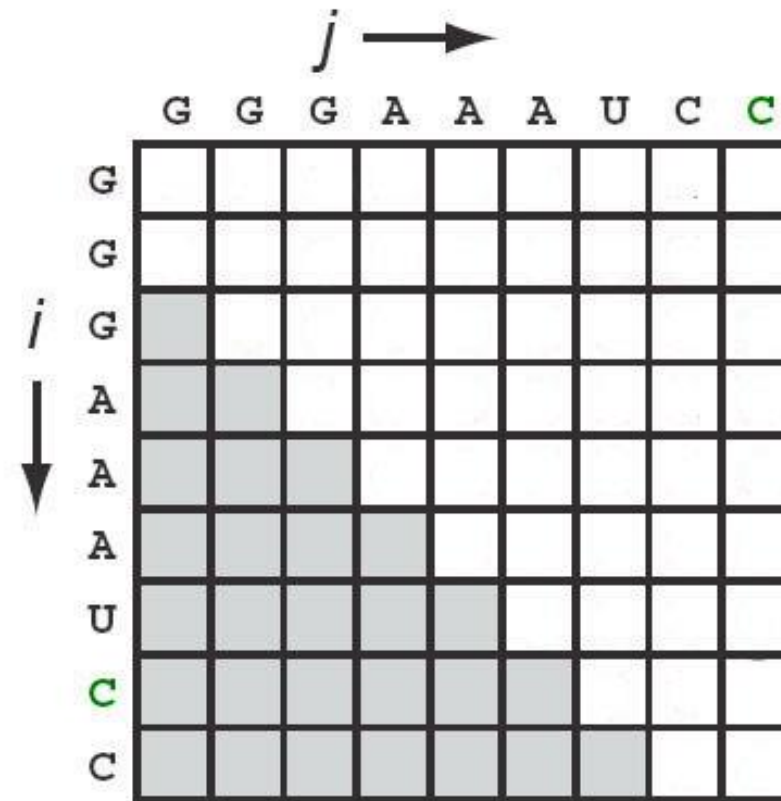
- Recurrence:

$$S(i, j) = \max \begin{cases} S(i+1, j-1) + 1 & \text{(if } (i, j) \text{ base pair)} \\ S(i+1, j) \\ S(i, j-1) \\ \max_{1 \leq k < j} S(i, k) + S(k+1, j) \end{cases}$$



Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension



Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

Initialize first two diagonals to 0

Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

Fill in squares sweeping diagonally

	G	G	G	A	A	A	U	C	C
G	0								
G	0	0							
G		0	0						
A			0	0					
A				0	0				
A					0	0			
U						0	0		
C							0	0	○
C								0	0

Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

Bases cannot pair

Diagram illustrating a sequence alignment matrix (likely edit distance) for nucleotide sequences. The columns are labeled j (horizontal axis) and the rows are labeled i (vertical axis). The matrix shows values for sequences G, G, G, A, A, A, U, C, C, C. The bottom-right cell (C, C) is highlighted with a circle, indicating the final state.

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
A	1	1	0	0	0	0	1	1	1
A	1	1	0	0	0	0	1	1	1
A	1	1	0	0	0	0	1	1	1
U	2	2	1	1	1	0	0	0	0
C	3	3	2	2	2	1	0	0	0
C	3	3	2	2	2	1	0	0	0

Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

Bases can pair, similar to matched alignment

Diagram illustrating a sequence alignment matrix (likely edit distance) for two sequences. The columns are labeled G, G, G, A, A, A, U, C, C, C and the rows are labeled G, G, G, A, A, A, U, C, C, C. The matrix is filled with values representing the edit distance. The top-left cell is 0. The values increase as the edit distance increases. The bottom-right cell is 9. The matrix is symmetric along the main diagonal.

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
A	1	1	1	0	0	0	1	1	2
A	1	1	1	0	0	0	1	1	2
A	1	1	1	0	0	0	1	1	2
U	2	2	2	1	1	1	0	0	1
C	3	3	3	2	2	2	1	0	0
C	3	3	3	2	2	2	1	0	0

Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

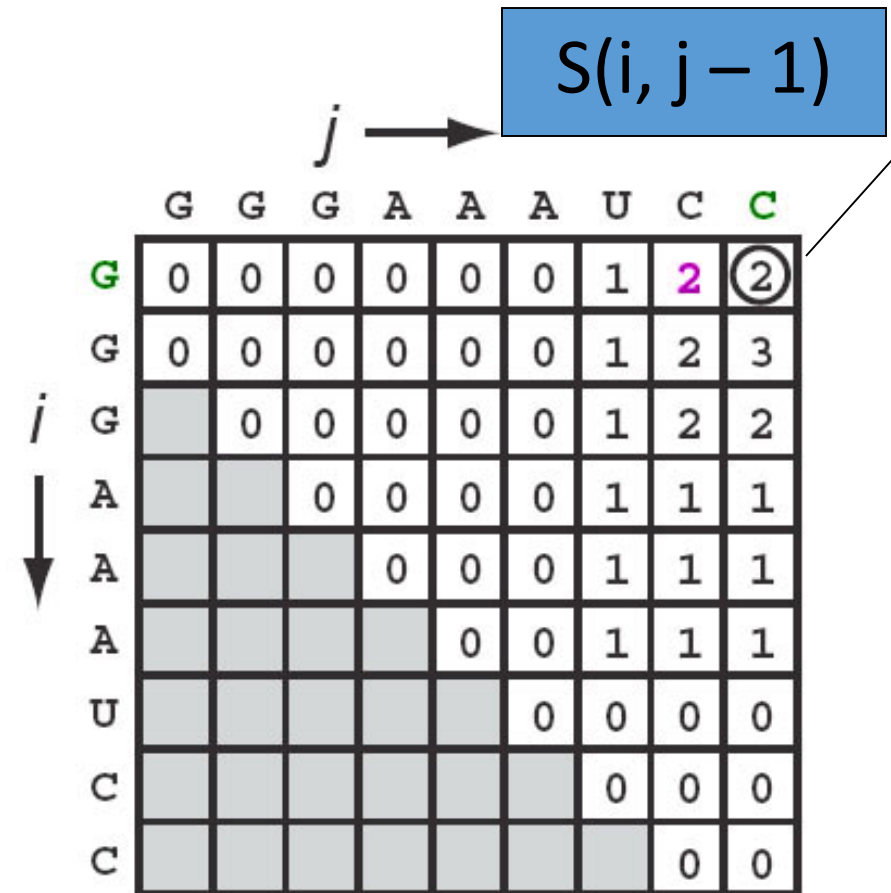
	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	②
G	0	0	0	0	0	0	1	2	3
G		0	0	0	0	0	1	2	2
A			0	0	0	0	1	1	1
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
C							0	0	0
C								0	0

Dynamic Programming—
possible paths

Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

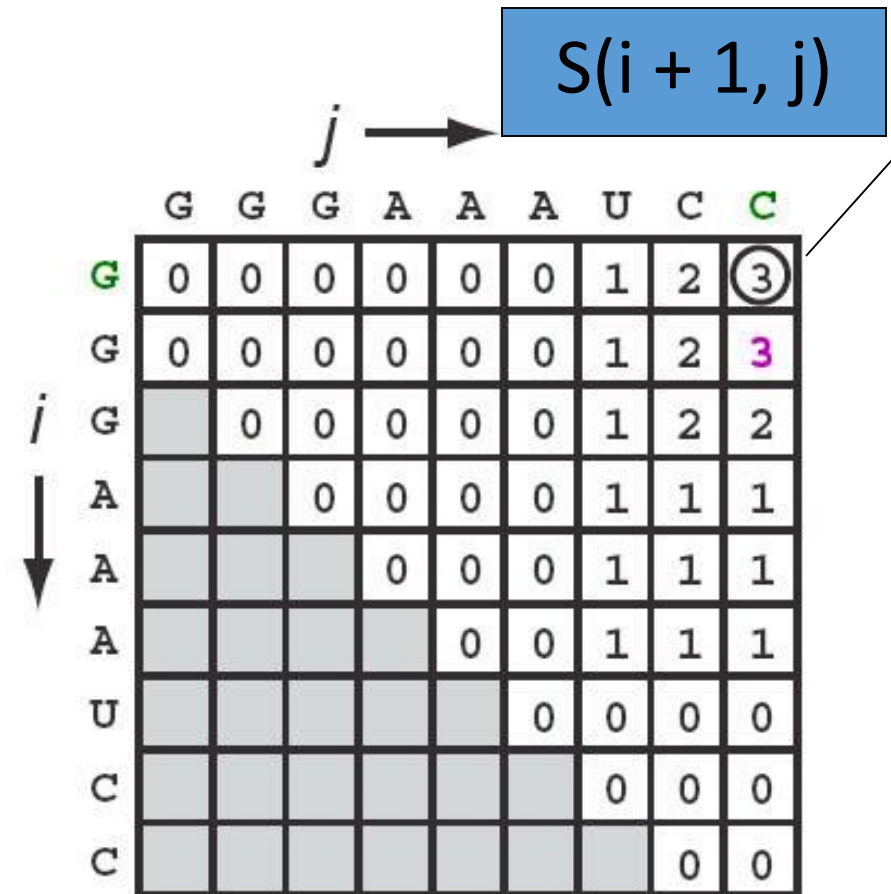
Dynamic Programming—
possible paths



Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

Dynamic Programming—
possible paths



Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

Dynamic Programming—
possible paths

Images—Sean Eddy

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	②
G	0	0	0	0	0	0	1	2	3
G		0	0	0	0	0	1	2	2
A			0	0	0	0	1	1	1
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
C							0	0	0
C								0	0

$$S(i + 1, j - 1) + 1$$

Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

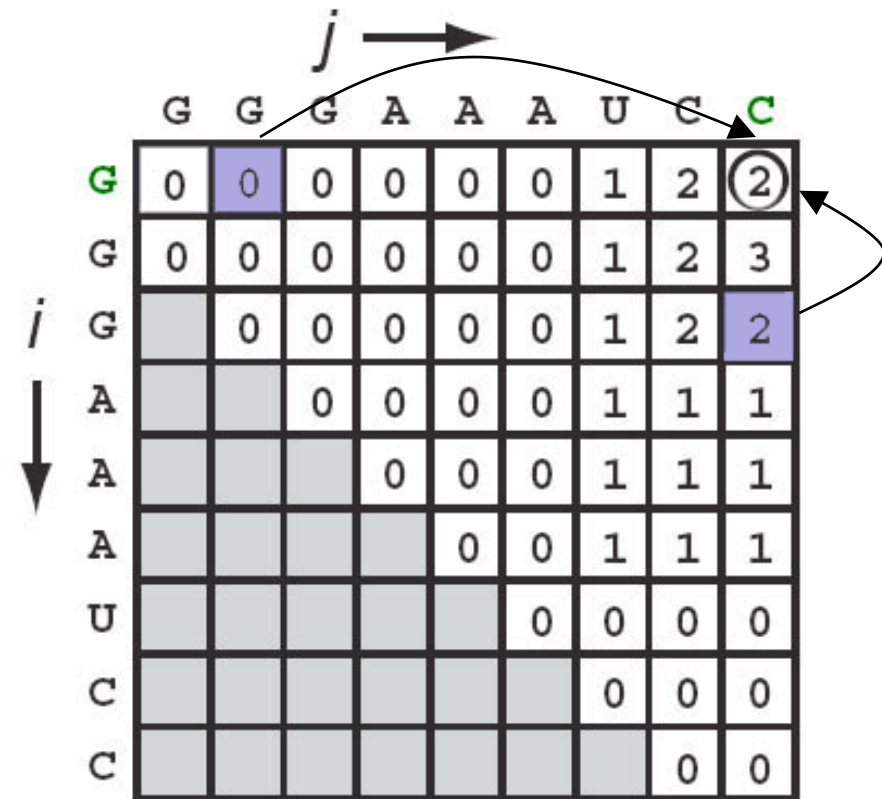
				$j \rightarrow$						
		G	G	G	A	A	A	U	C	C
$i \downarrow$	G	0	0	0	0	0	0	1	2	3
	G	0	0	0	0	0	0	1	2	3
	G		0	0	0	0	0	1	2	2
	A			0	0	0	0	1	1	1
	A				0	0	0	1	1	1
	A					0	0	1	1	1
	U						0	0	0	0
	C							0	0	0
	C								0	0

Bifurcation—add values for all k

Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

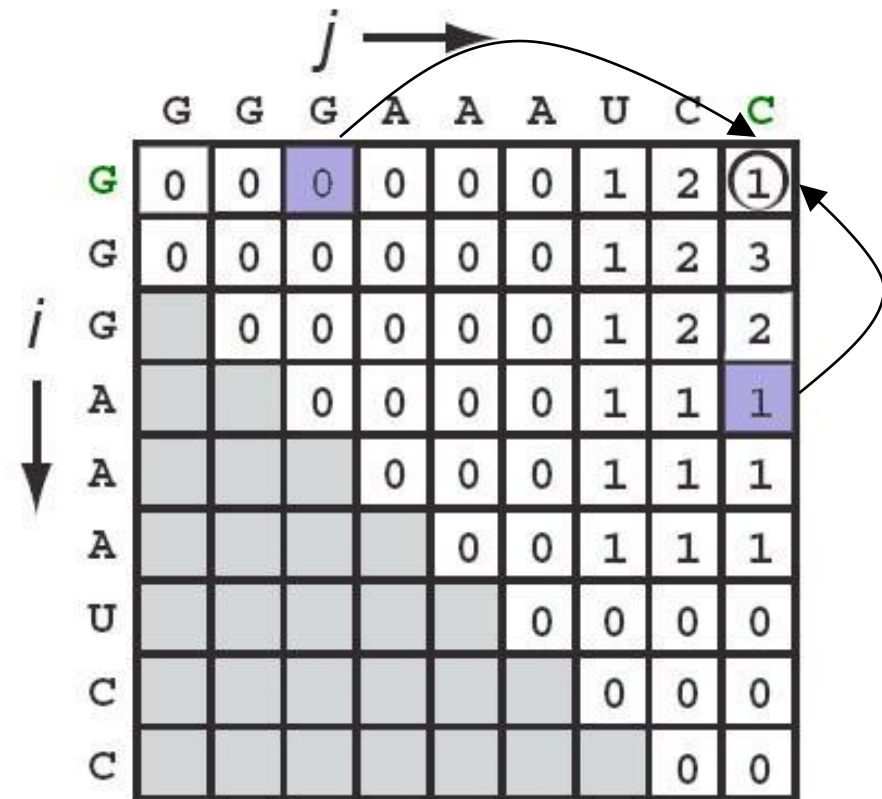
Bifurcation—add values for all k



Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

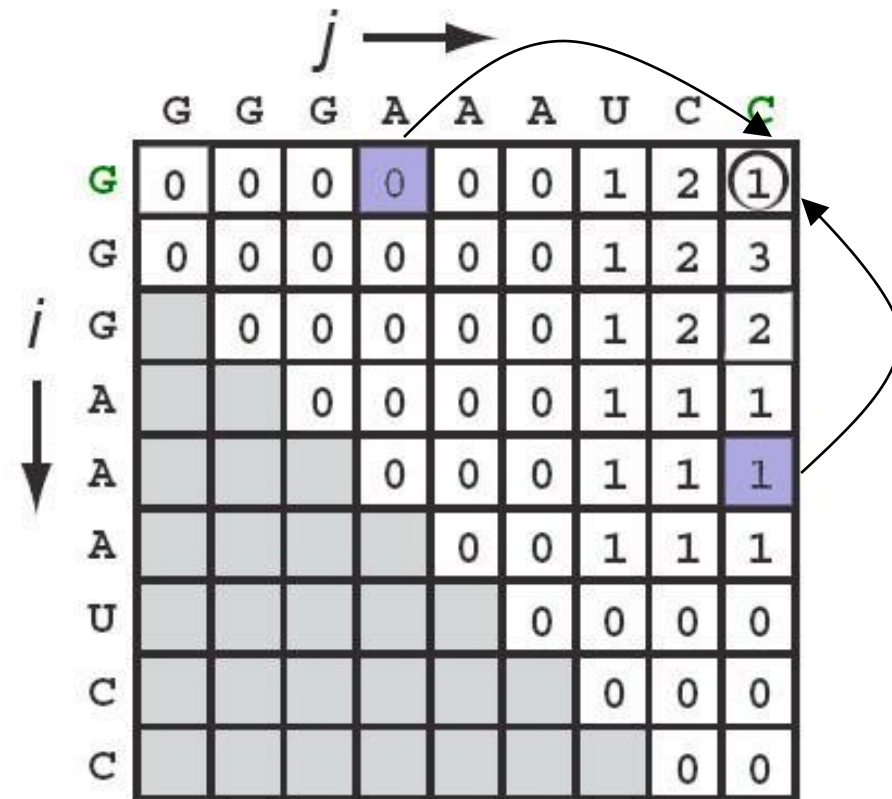
Bifurcation—add values for all k



Base Pair Maximization: Dynamic Programming

- Alignment Method:
 - Align RNA strand to itself
 - Score increases for feasible base pairs
- Each score independent of overall structure
- Bifurcation adds extra dimension

Bifurcation—add values for all k

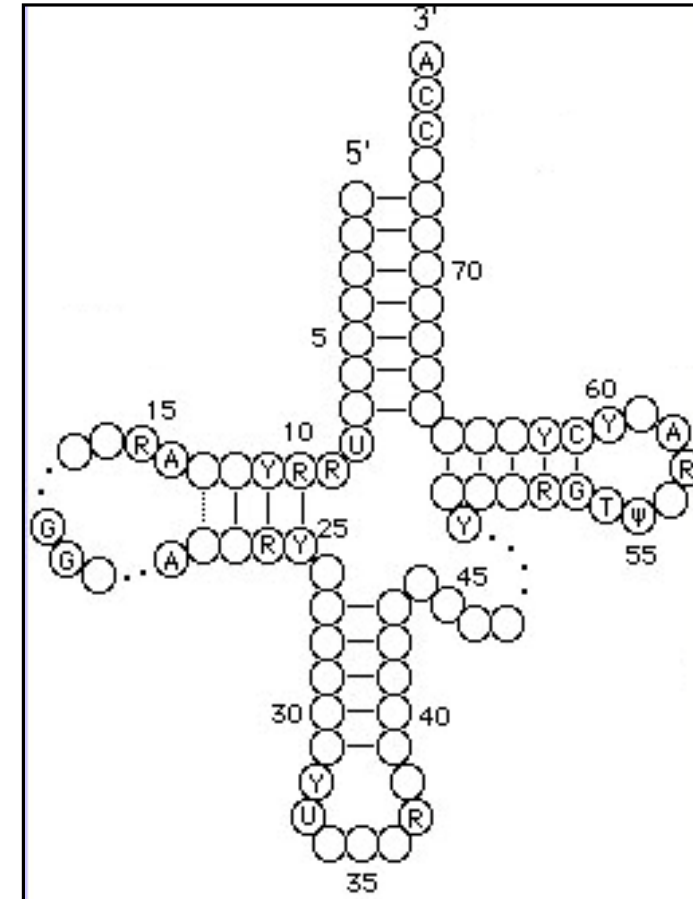


Base Pair Maximization: Drawbacks

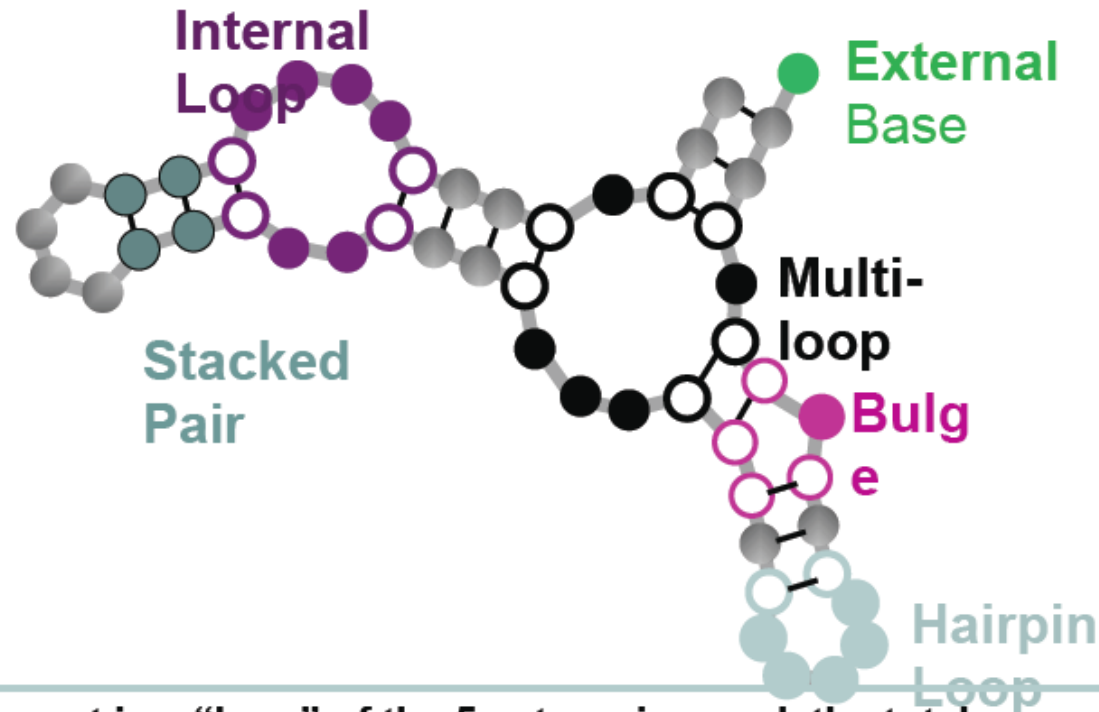
- Base pair maximization will not necessarily lead to the most stable structure.
 - It may create structure with many interior loops or hairpins which are energetically unfavorable.
- Results comparable to aligning sequences with scattered matches—not biologically reasonable.

Energy Minimization

- Thermodynamic Stability
 - Estimated using experimental techniques.
 - Theory : Most Stable = Most likely
- No pseudoknots due to algorithm limitations.
- Attempts to maximize the score, taking thermodynamics into account.
- MFOLD and ViennaRNA



Energy minimization for RNA prediction



Every element is a “loop” of the 5 categories. and the total energy is the sum of all loop energies



energy is assigned to substructures, not to base pairs

How to find the minimum energy secondary structure? (Zucker Algo)



Similar to Nussinov, we use DP

$W(j)$: energy of the optimal secondary structure for $S[1..j]$

$V(i, j)$: optimal energy for $S[i..j]$ with (i, j) forming base pair

$VBI(i, j)$: optimal energy for $S[i..j]$ with (i, j) closing an internal loop

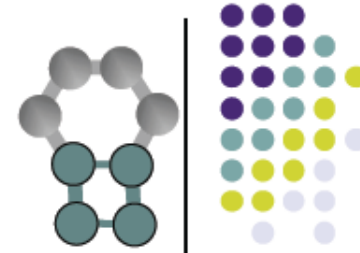
$VM(i, j)$: optimal energy for $S[i..j]$ with (i, j) closing a multi-loop



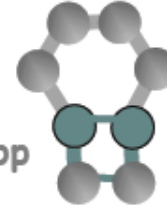
Loop energy

- $eS(i, j)$: free energy of the **stacking pair** consists of base pairs (i, j) and $(i+1, j-1)$. Stacking pair stabilizes the structure and has a negative energy
- eS depends on all the four bases involved in the stack
- $eH(i, j)$: free energy of the **hairpin** closed by base pair i, j
- Depends on loopsize; the unpaired bases adjacent to i, j in loop
- $eL(i, j, i', j')$: free energy of an **internal loop** or bulge enclosed by (i, j) and consists of 2 base pairs.
- Similar to eH ; depends on four paired bases (i, j, i', j') ; loopsize;
- Unpaired bases next to the paired bases
- $eM(i, j, i_1, j_1, \dots, i_k, j_k)$: free energy of a **multi-loop** enclosed by (i, j) and consists of $k+1$ base pairs.
- Similar to eL , but many approximations used.

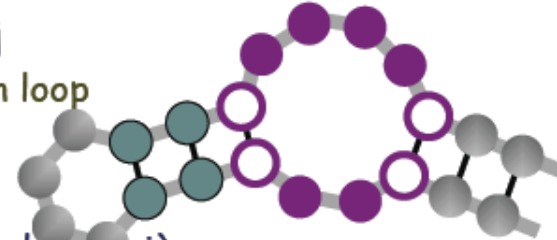
Stacked Pair

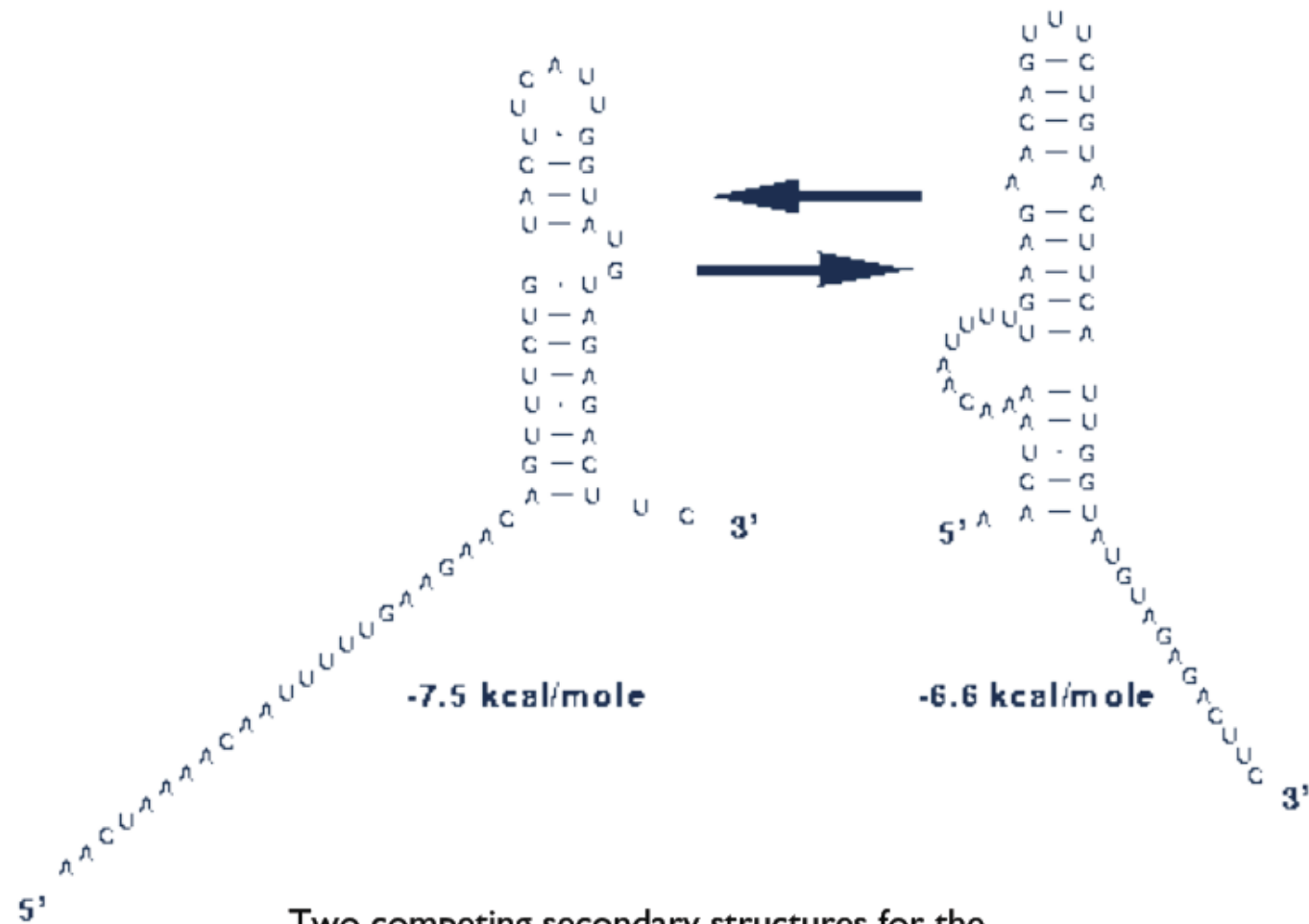


Hairpin loop



Internal Loop

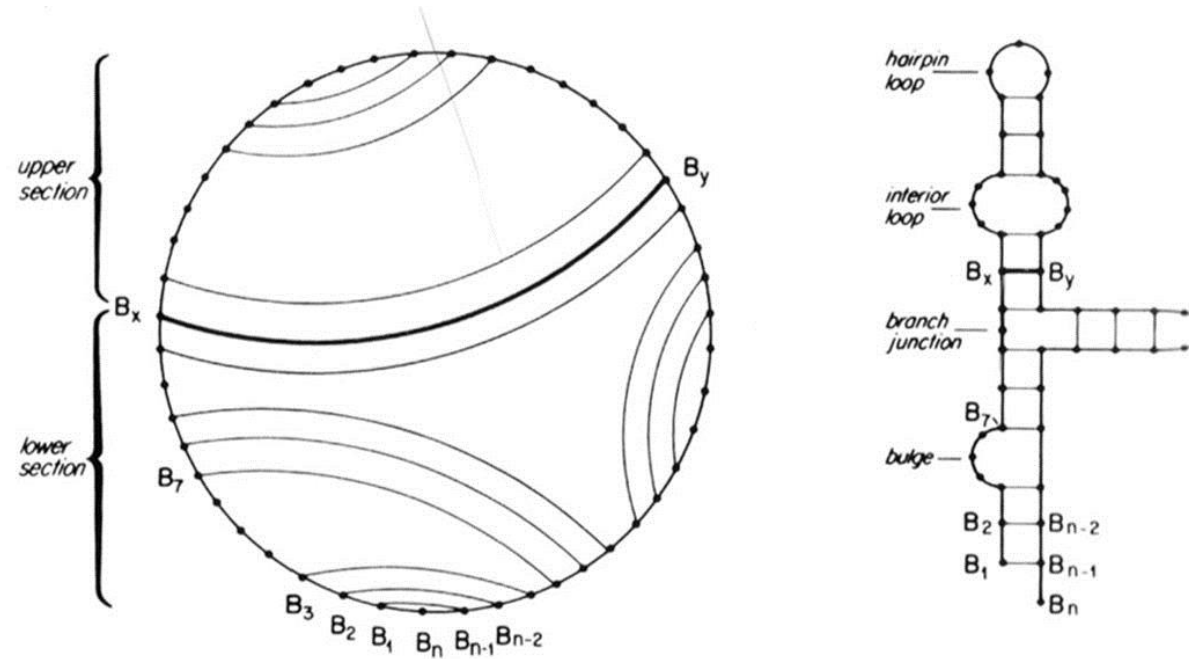




Two competing secondary structures for the *Leptomonas collosoma* spliced leader mRNA.

Energy Minimization Results

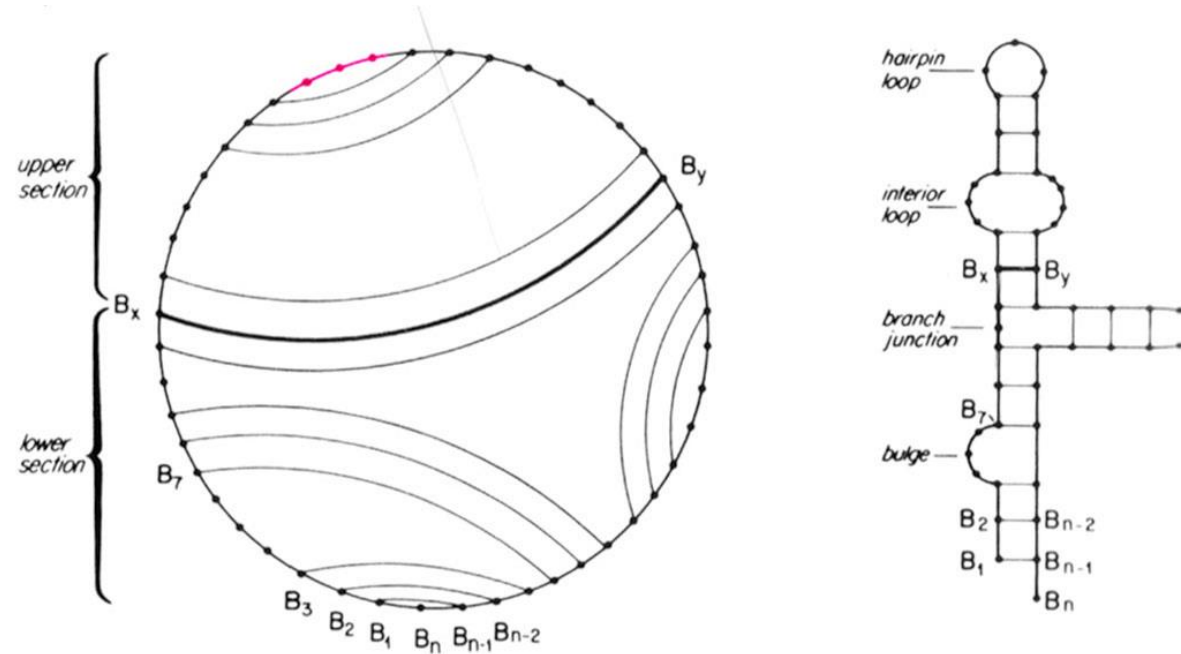
- Linear RNA strand folded back on itself to create secondary structure
- Circularized representation uses this requirement
 - Arcs represent base pairing



Images – David Mount

Energy Minimization Results

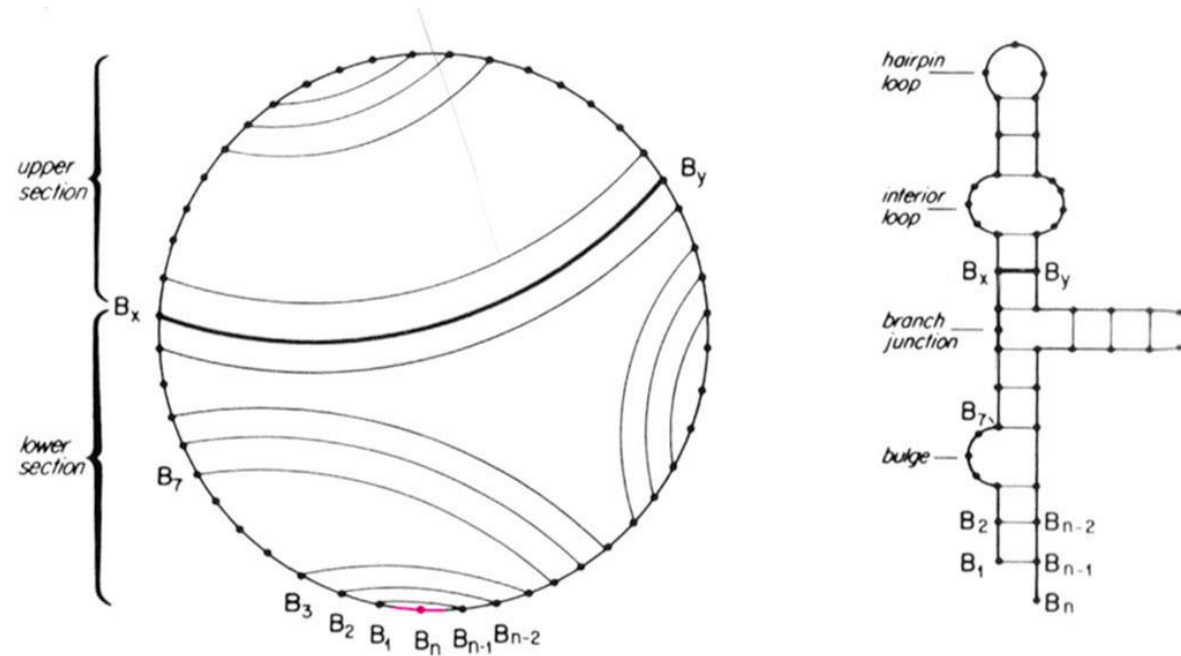
- All loops must have exactly three bases in them.
 - Equivalent to having at least three base pairs between arc endpoints.



Images – David Mount

Energy Minimization Results

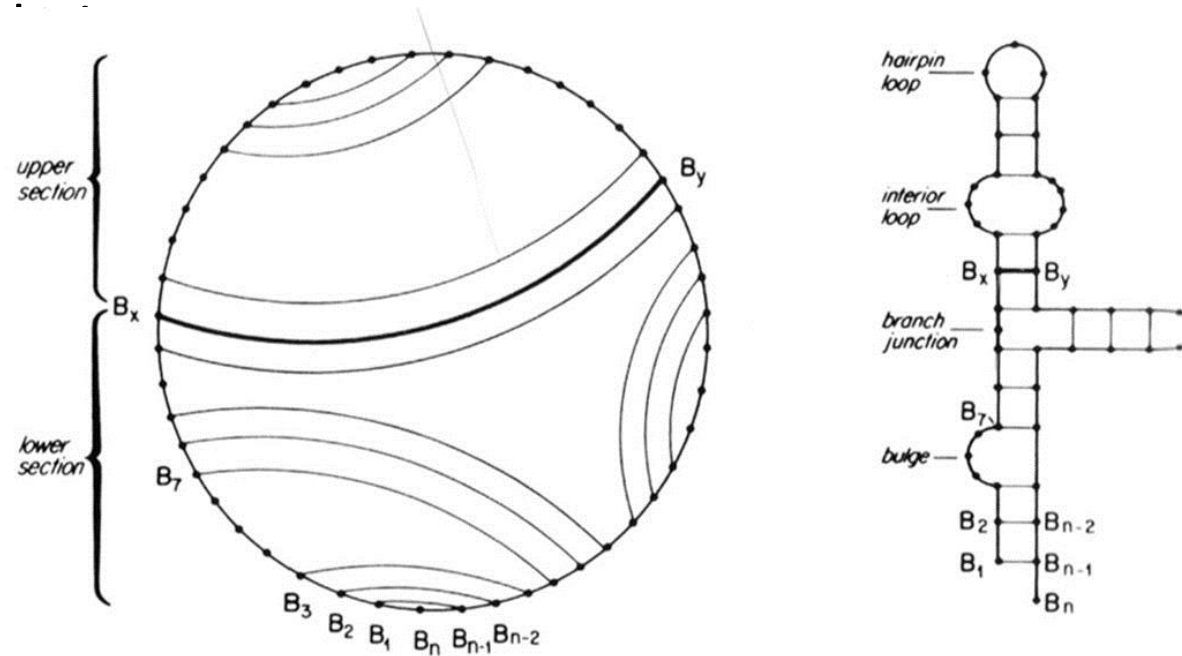
- All loops must have exactly three bases in them.
 - **Exception:** Location where beginning and end of RNA come together in circularized representation.



Images – David Mount

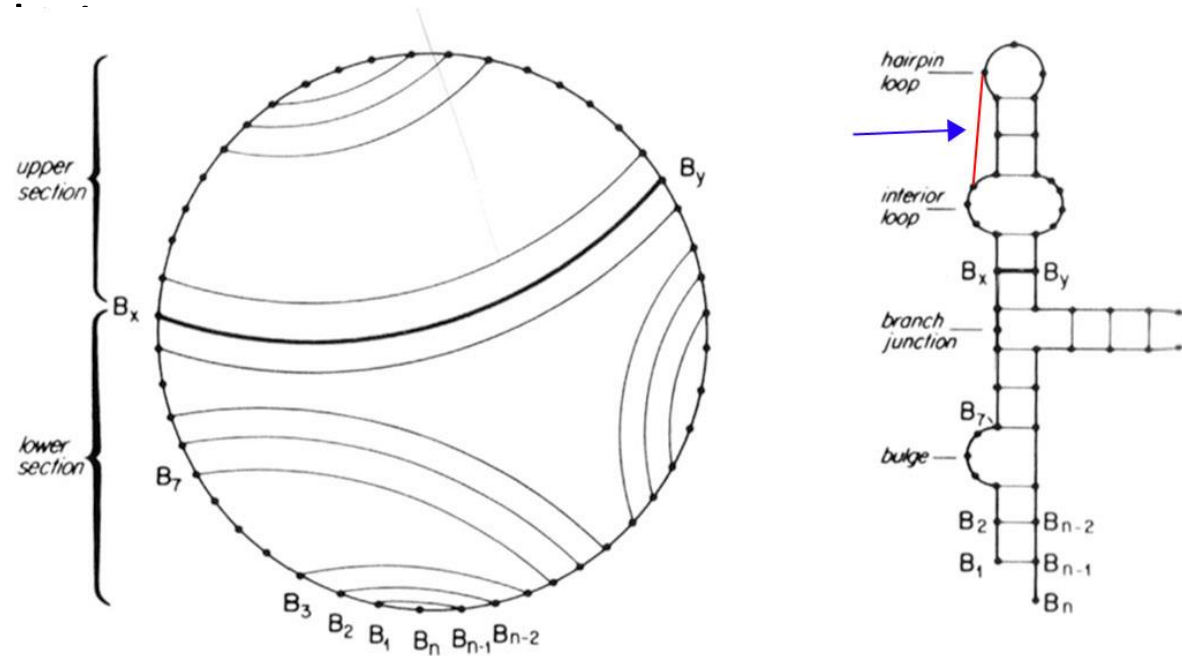
Trouble with Pseudoknots

- Pseudoknots cause a breakdown in the dynamic programming algorithm.
- In order to form a pseudoknot, checks must be made to ensure base is not already paired—this breaks down the recurrence



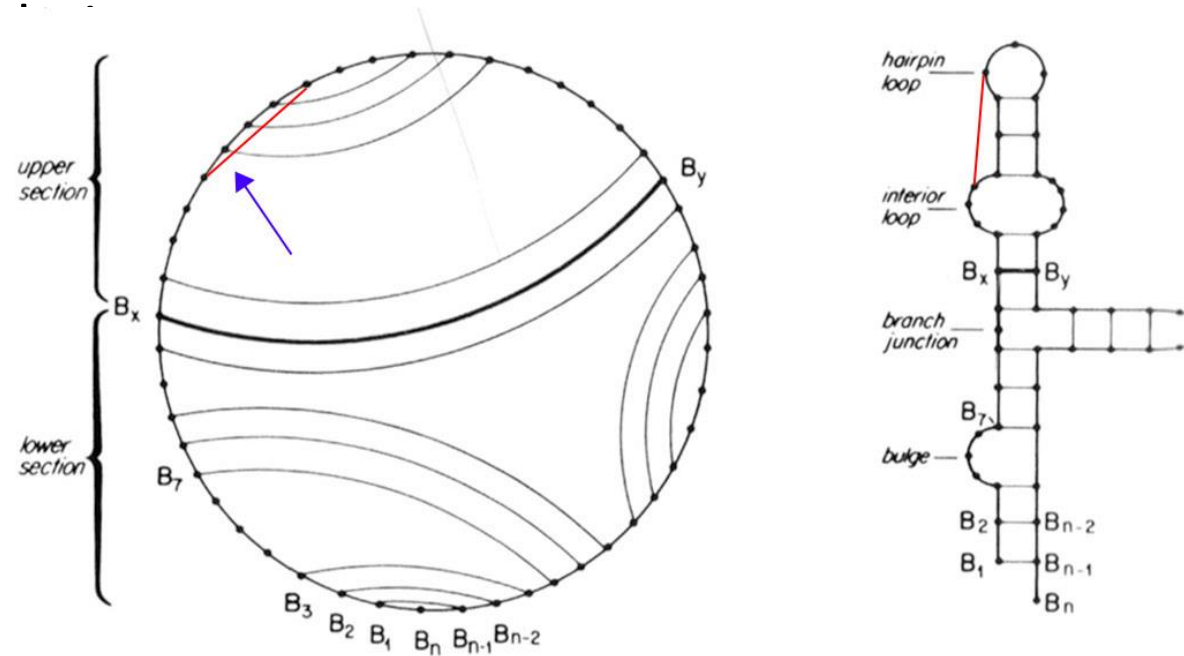
Trouble with Pseudoknots

- Pseudoknots cause a breakdown in the dynamic programming algorithm.
- In order to form a pseudoknot, checks must be made to ensure base is not already paired—this breaks down the recurrence



Trouble with Pseudoknots

- Pseudoknots cause a breakdown in the dynamic programming algorithm.
- In order to form a pseudoknot, checks must be made to ensure base is not already paired—this breaks down the recurrence

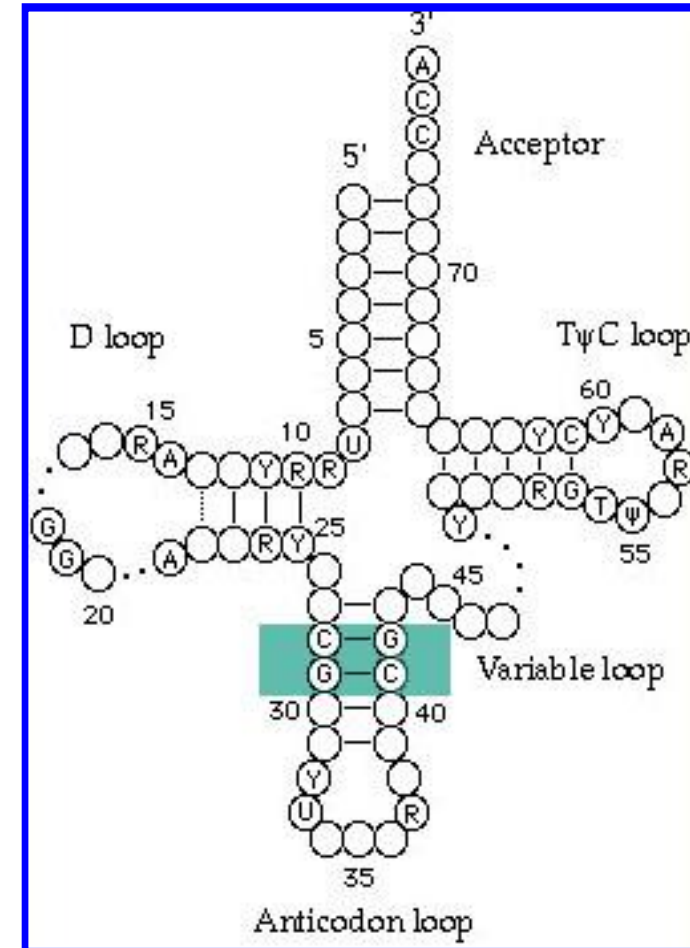


Energy Minimization: Drawbacks

- Computes only one optimal structure.
- Optimal solution may not represent the biologically correct solution.

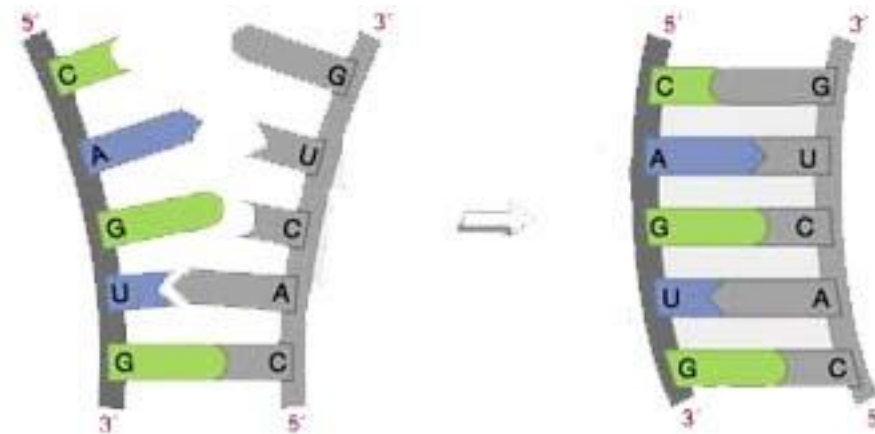
Alternative Algorithms - Covariation

- Expect areas of base pairing in tRNA to be covarying between various species. Find this by sequence alignment.
- Base pairing creates same stable tRNA structure in organisms.



Sequence Alignment to Determine Structure

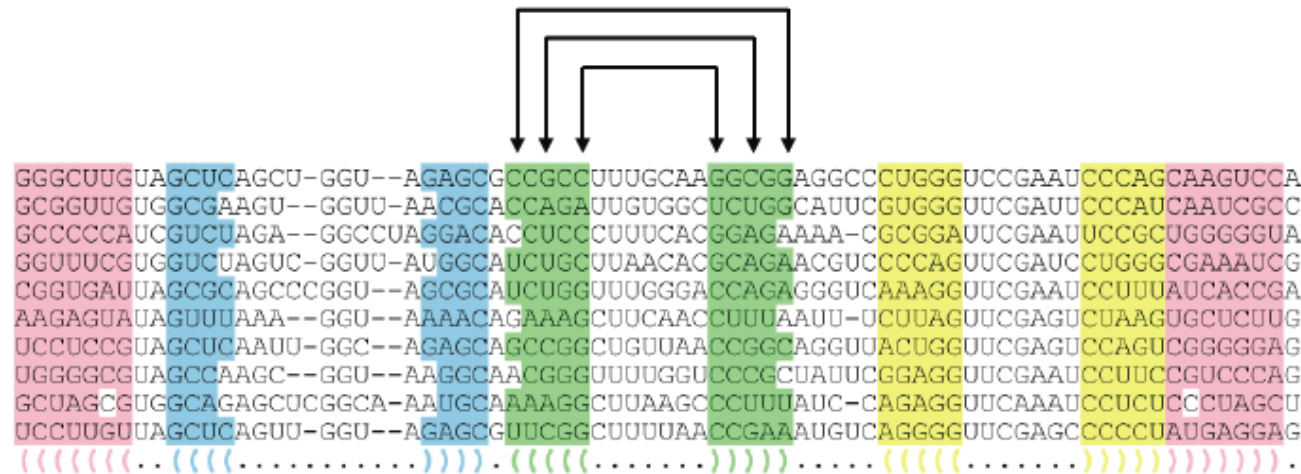
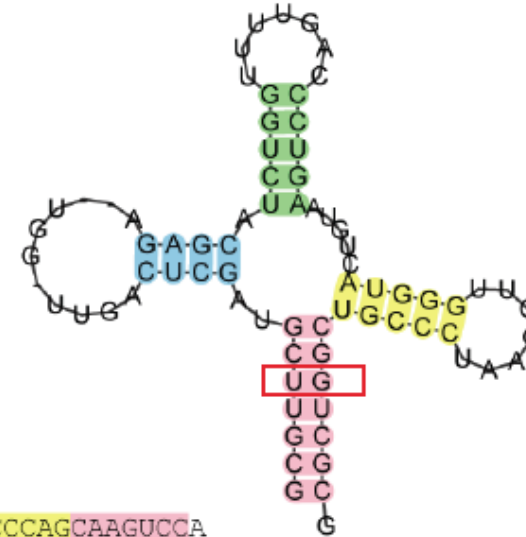
- Bases pair in order to form backbones and determine the secondary structure.
- Aligning bases based on their ability to pair with each other gives an algorithmic approach to determining the optimal structure.



Mutual Information in RNA structure prediction



$$I(X,Y) = \sum_i \sum_j f_{X,Y}(x_i, y_j) \cdot \log \frac{f_{X,Y}(x_i, y_j)}{f_X(x_i) \cdot f_Y(y_j)}$$



IV

	1	2	3	4	5	6	7	8	9
A	G	A	U	A	A	U	C	U	
A	G	A	U	C	A	U	C	U	
A	G	A	C	G	U	U	C	U	
A	G	A	U	U	U	U	C	U	
A	G	C	C	A	G	G	C	U	
A	G	C	G	C	G	G	C	U	
A	G	C	U	G	C	G	C	U	
A	G	C	A	U	C	G	C	U	
A	G	G	U	A	G	C	C	U	
A	G	G	G	C	G	C	C	U	
A	G	G	U	G	U	C	C	U	
A	G	G	C	U	U	C	C	U	
A	G	U	A	A	A	A	C	U	
A	G	U	C	C	A	A	C	U	
A	G	U	U	G	C	A	C	U	
A	G	U	U	U	C	A	C	U	

	1	2	3	4	5	6	7	8	9
A	16	0	4	2	4	4	4	0	0
C	0	0	4	4	4	4	4	16	0
G	0	16	4	2	4	4	4	0	0
U	0	0	4	8	4	4	4	0	16

[illegible]

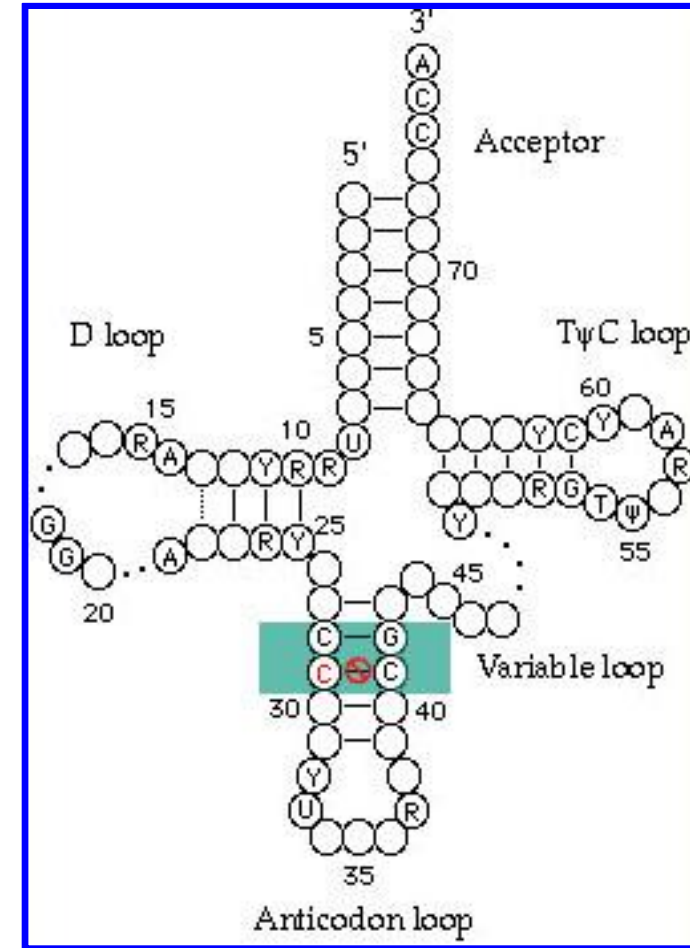
Cols 1 & 9, 2 & 8: perfect conservation and *might* be base-paired, but unclear whether they are. M.I. = 0

Cols 3 & 7: completely unconserved, but always W-C pairs, so seems likely that they do base-pair. M.I. = 2 bits.

Cols 7->6: unconserved, but each letter in 7 has only 2 possible mates in 6. M.I. = 1 bit.

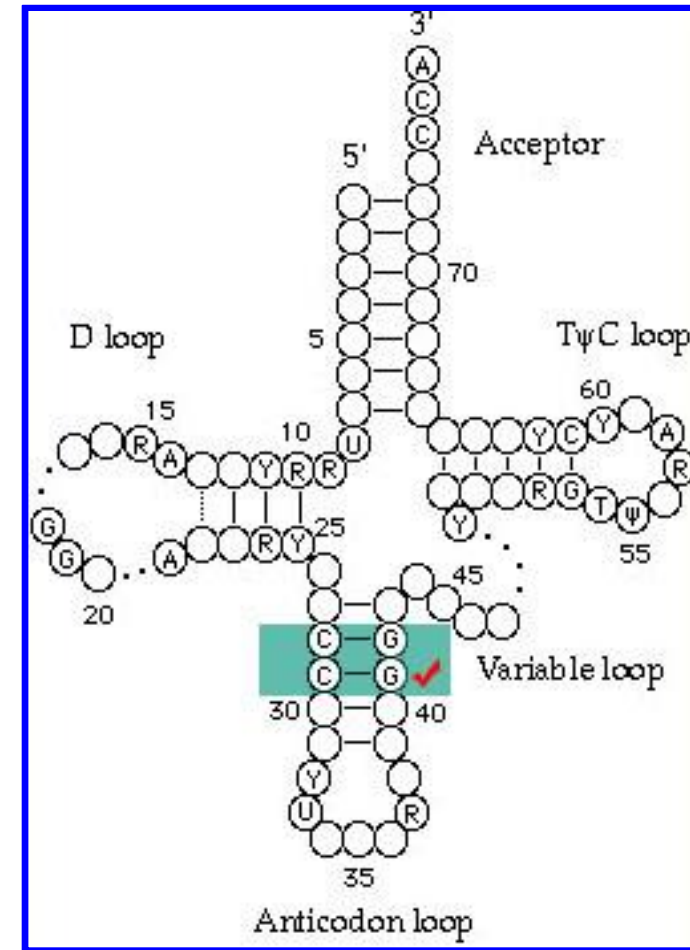
Alternative Algorithms - Covariation

- Expect areas of base pairing in tRNA to be covarying between various species.
- Base pairing creates same stable tRNA structure in organisms.
- Mutation in one base yields pairing impossible and breaks down structure.



Alternative Algorithms - Covariation

- Expect areas of base pairing in tRNA to be covarying between various species.
- Base pairing creates same stable tRNA structure in organisms.
- Mutation in one base yields pairing impossible and breaks down structure.
- Covariation ensures ability to base pair is maintained and RNA structure is conserved.



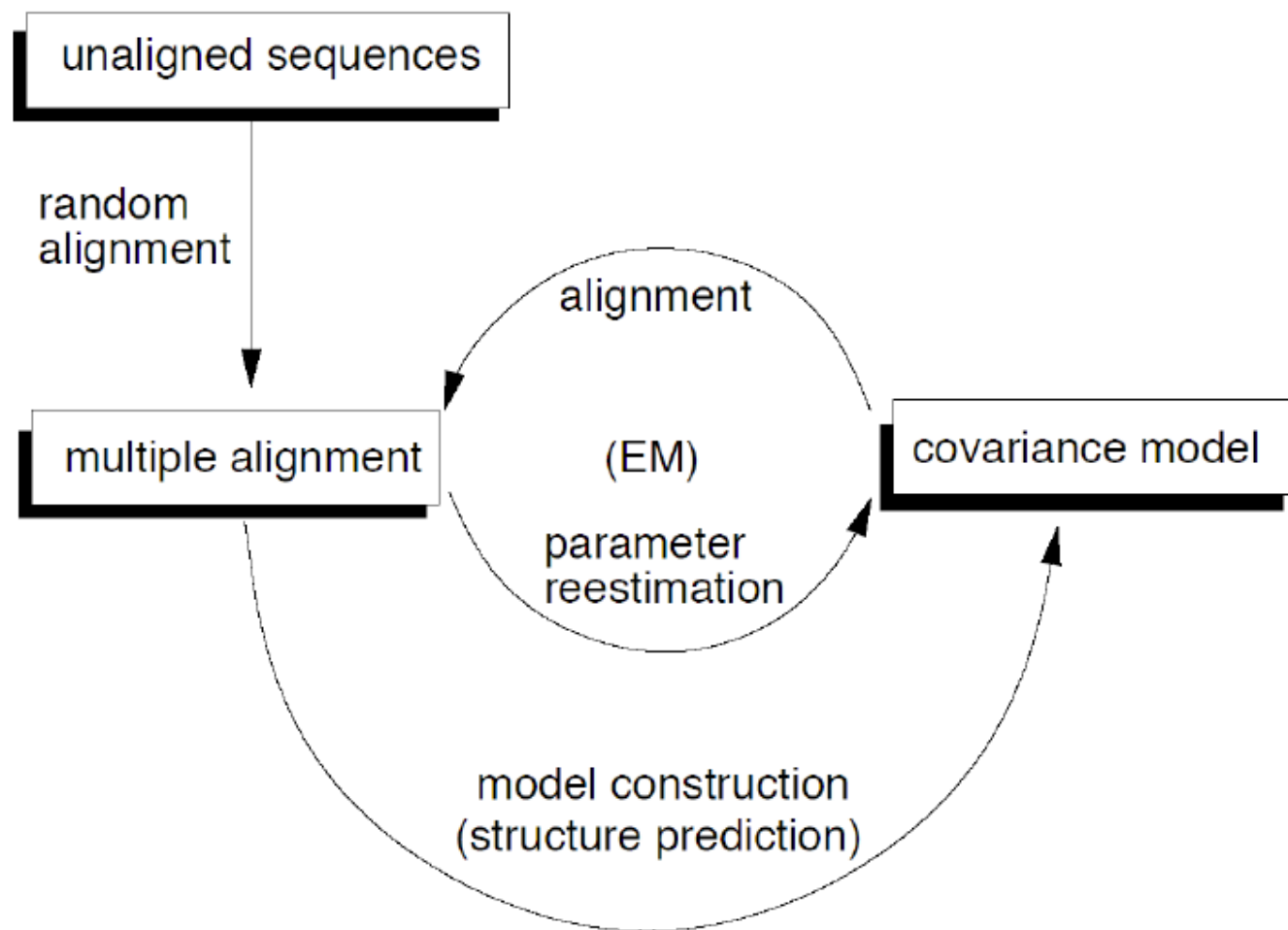
MI-Based Structure-Learning

- find best (max total MI) subset of column pairs among $i \dots j$, subject to absence of pseudo-knots

$$S_{i,j} = \max \begin{cases} S_{i+1,j} \\ S_{i,j-1} \\ S_{i+1,j-1} + M_{i,j} \\ \max_{i < j < k} S_{i,k} + S_{k+1,j} \end{cases}$$

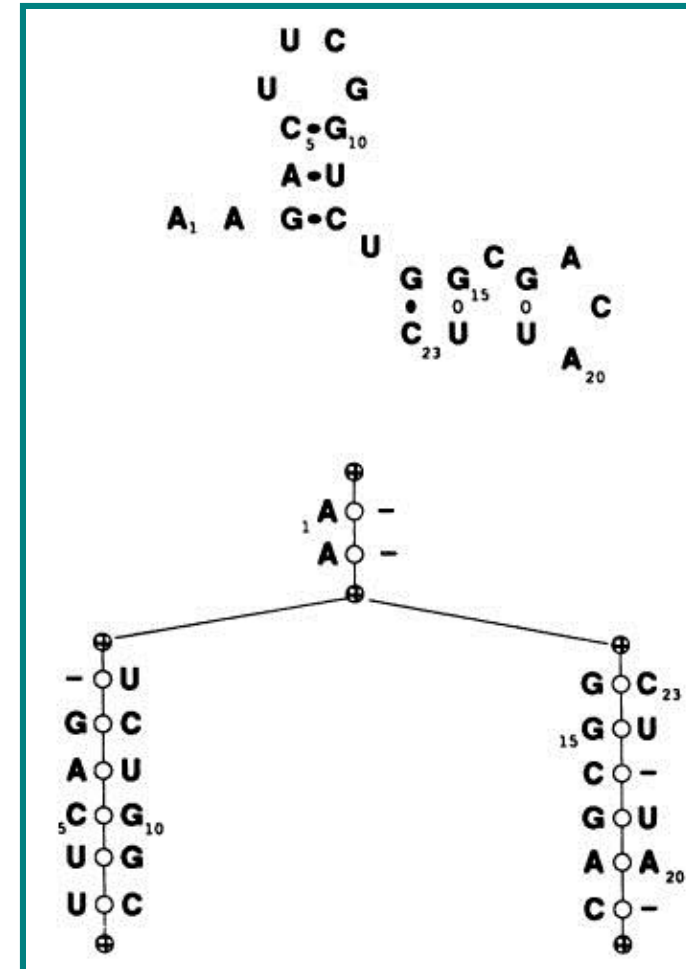
- “just like Nussinov/Zucker folding”
- BUT, need enough data---enough sequences at right phylogenetic distance

Model Training



Binary Tree Representation of RNA Secondary Structure

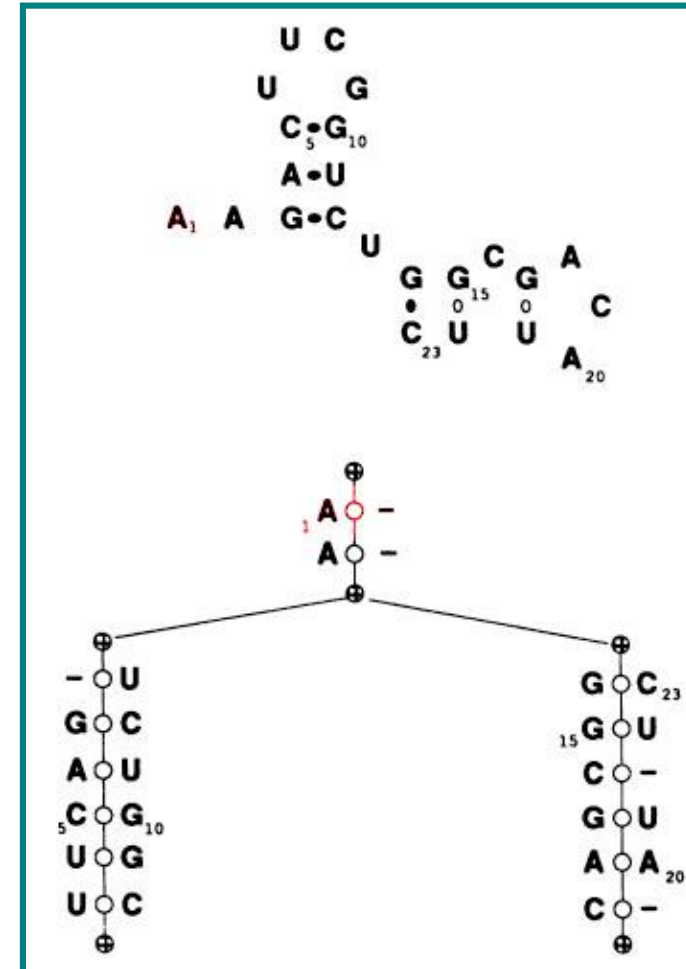
- Representation of RNA structure using Binary tree
- Nodes represent
 - Base pair if two bases are shown
 - Loop if base and “gap” (dash) are shown
- Pseudoknots still not represented
- Tree does not permit varying sequences
 - Mismatches
 - Insertions & Deletions



Images – Eddy et al.

Binary Tree Representation of RNA Secondary Structure

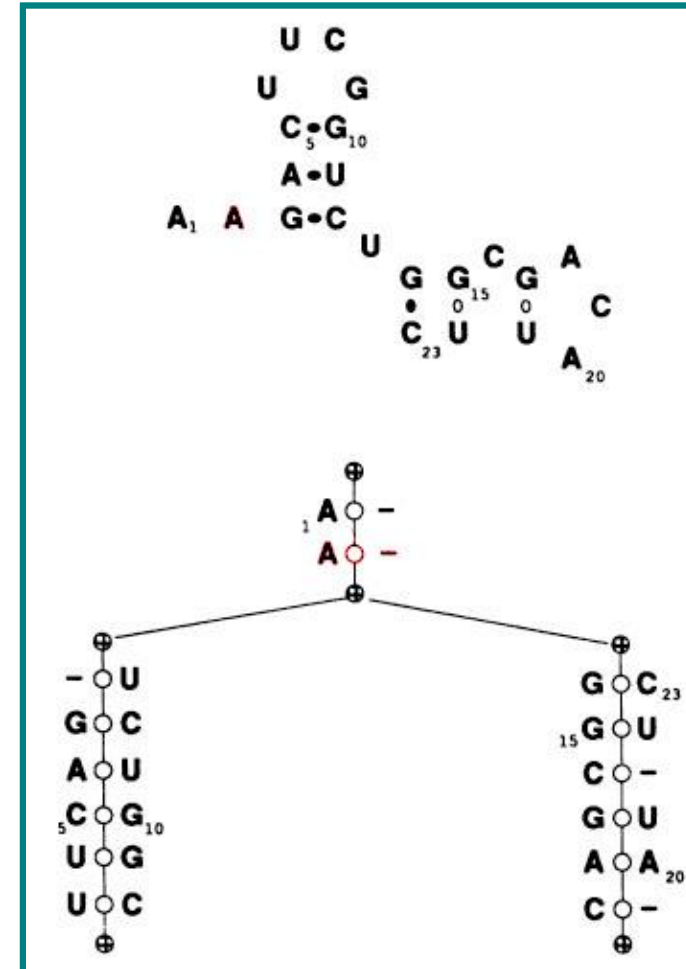
- Representation of RNA structure using Binary tree
- Nodes represent
 - Base pair if two bases are shown
 - Loop if base and “gap” (dash) are shown
- Pseudoknots still not represented
- Tree does not permit varying sequences
 - Mismatches
 - Insertions & Deletions



Images – Eddy et al.

Binary Tree Representation of RNA Secondary Structure

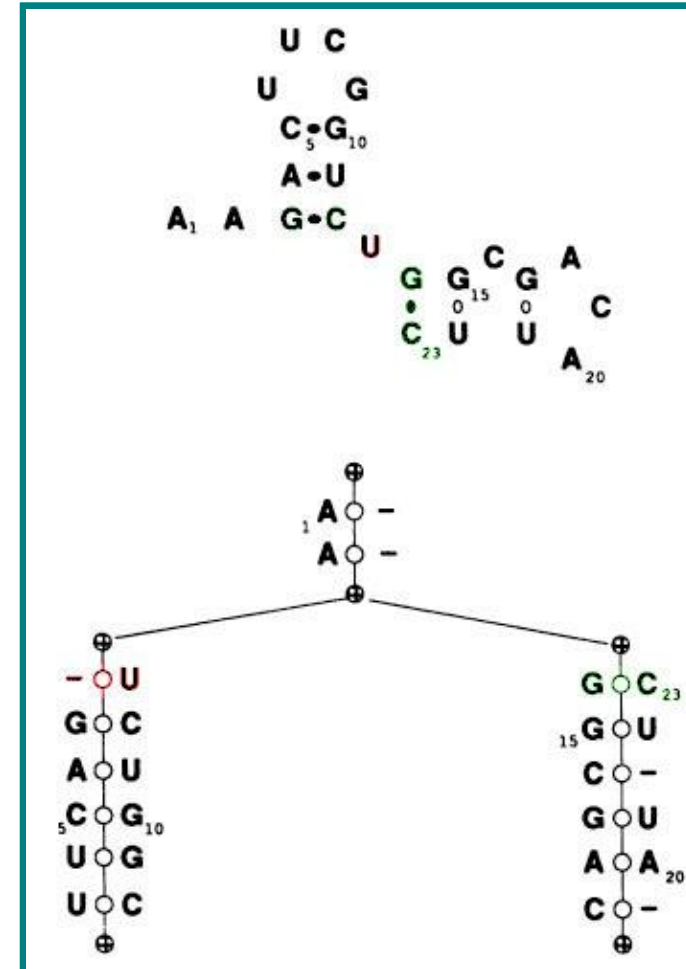
- Representation of RNA structure using Binary tree
- Nodes represent
 - Base pair if two bases are shown
 - Loop if base and “gap” (dash) are shown
- Pseudoknots still not represented
- Tree does not permit varying sequences
 - Mismatches
 - Insertions & Deletions



Images – Eddy et al.

Binary Tree Representation of RNA Secondary Structure

- Representation of RNA structure using Binary tree
- Nodes represent
 - Base pair if two bases are shown
 - Loop if base and “gap” (dash) are shown
- Pseudoknots still not represented
- Tree does not permit varying sequences
 - Mismatches
 - Insertions & Deletions



Images – Eddy et al.

Covariance Model

- **Covariance Model:** HMM which permits flexible alignment to an RNA structure – emission and transition probabilities
- Model trees based on finite number of states
 - Match states – sequence conforms to the model:
 - MATP: State in which bases are paired in the model and sequence.
 - MATL & MATR: State in which either right or left bulges in the sequence and the model.
 - Deletion – State in which there is deletion in the sequence when compared to the model.
 - Insertion – State in which there is an insertion relative to model.

Covariance Model

- **Covariance Model:** HMM which permits flexible alignment to an RNA structure – emission and transition probabilities
- Transitions have probabilities.
 - Varying probability: Enter insertion, remain in current state, etc.
 - Bifurcation: No probability, describes path.

Covariance Model (CM) Training Algorithm

- $S(i, j)$ = Score at indices i and j in RNA when aligned to the Covariance Model.

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

$$S(i, j) = \max \left\{ \begin{array}{l} S(i+1, j-1) + M(i, j) \\ S(i+1, j) \\ S(i, j-1) \\ \max_{i < k < j} S(i, k) + S(k+1, j) \end{array} \right.$$

- Frequencies obtained by aligning model to “training data”—consists of sample sequences.
 - Reflect values which optimize alignment of sequences to model.

Covariance Model (CM) Training Algorithm

- $S(i, j)$ = Score at indices i and j in RNA when aligned to the Covariance Model.

$$S(i, j) = \max \begin{cases} S(i+1, j-1) + M(i, j) \\ S(i+1, j) \\ S(i, j-1) \\ \max_{i < k < j} S(i, k) + S(k+1, j) \end{cases}$$

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

Frequency of seeing the symbols
(A, C, G, U) together in locations i and j
depending on symbol.

- Frequencies obtained by aligning model to “training data”—consists of sample sequences.
 - Reflect values which optimize alignment of sequences to model.

Covariance Model (CM) Training Algorithm

- $S(i, j)$ = Score at indices i and j in RNA when aligned to the Covariance Model.

$$S(i, j) = \max \begin{cases} S(i+1, j-1) + M(i, j) \\ S(i+1, j) \\ S(i, j-1) \\ \max_{i < k < j} S(i, k) + S(k+1, j) \end{cases}$$

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

Independent frequency of seeing the symbols (A, C, G, T) in locations i or j depending on symbol.

- Frequencies obtained by aligning model to “training data”—consists of sample sequences.
 - Reflect values which optimize alignment of sequences to model.

Covariance Model (CM) Training Algorithm

- $S(i, j)$ = Score at indices i and j in RNA when aligned to the Covariance Model.

$$S(i, j) = \max \begin{cases} S(i+1, j-1) + M(i, j) \\ S(i+1, j) \\ S(i, j-1) \\ \max_{i < k < j} S(i, k) + S(k+1, j) \end{cases}$$

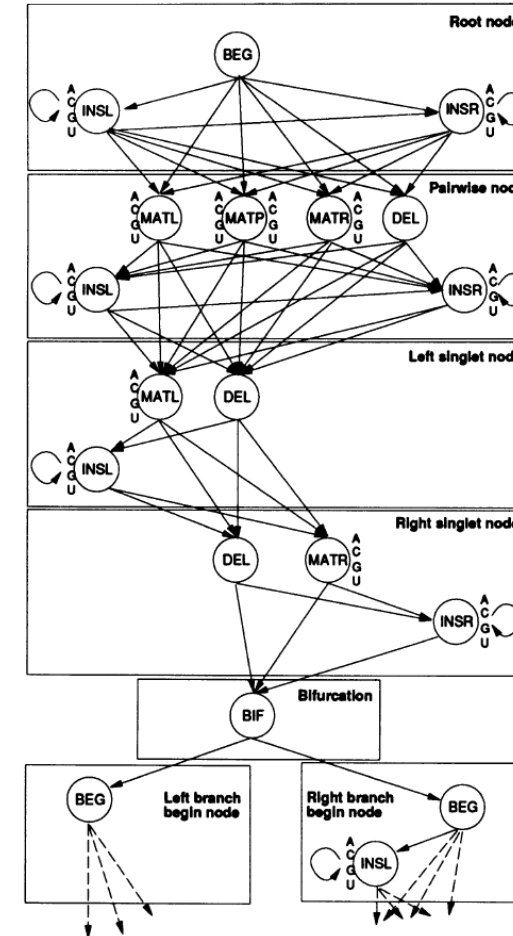
$$M_{i,j} = \sum_{x_i, x_j} f_{x_i x_j} \log_2 \frac{f_{x_i x_j}}{f_{x_i} f_{x_j}}$$

Independent frequency of seeing the symbols (A, C, G, U) in locations i or j depending on symbol.

- Frequencies obtained by aligning model to “training data”—consists of sample sequences.
 - Reflect values which optimize alignment of sequences to model.

Alignment to CM Algorithm

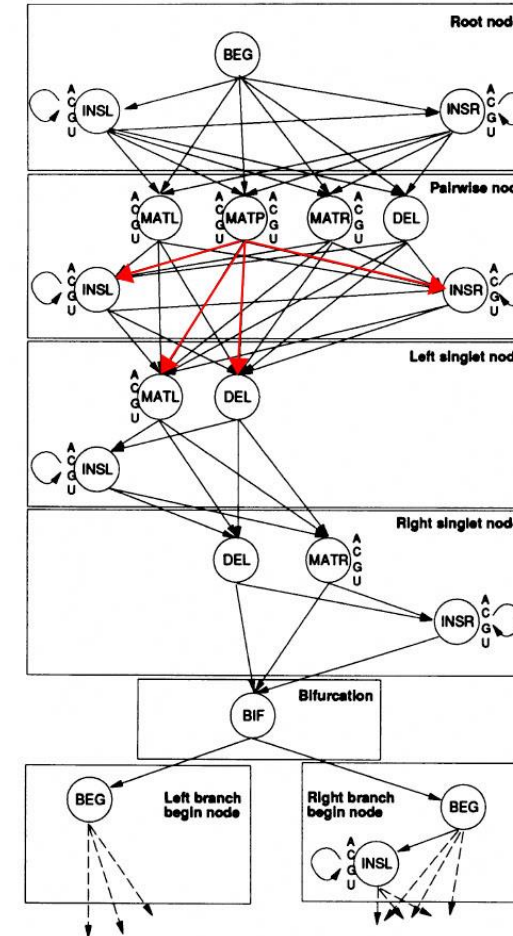
- Calculate the probability score of aligning RNA to CM.
- Three dimensional matrix— $O(n^3)$
 - Align sequence to given subtrees in CM.
 - For each subsequence, calculate all possible states.
- Subtrees evolve from bifurcations
 - For simplicity, left singlet is default.



Images—Eddy et al.

Alignment to CM Algorithm

- For each calculation, take into account:
 - Transition (T) to next state.
 - Emission probability (P) in the state as determined by training data.

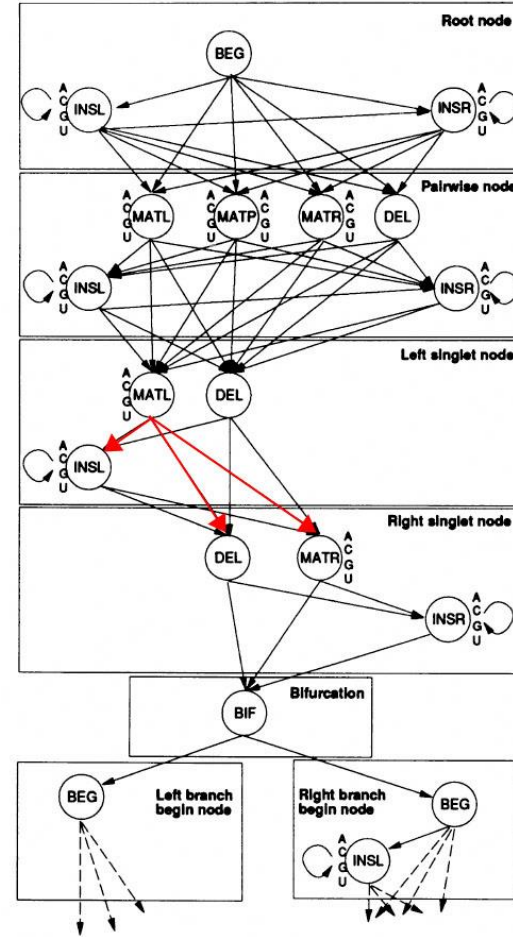


$$S_{i,j,y}(y = MATP) = \max_{y_{next}} [S_{i+1,j-1,y_{next}} + \log T(y_{next} | y) + \log \mathcal{P}(x_i, x_j | y)]$$

Images—Eddy et al.

Alignment to CM Algorithm

- For each calculation, take into account:
 - Transition (T) to next state.
 - Emission probability (P) in the state as determined by training data.

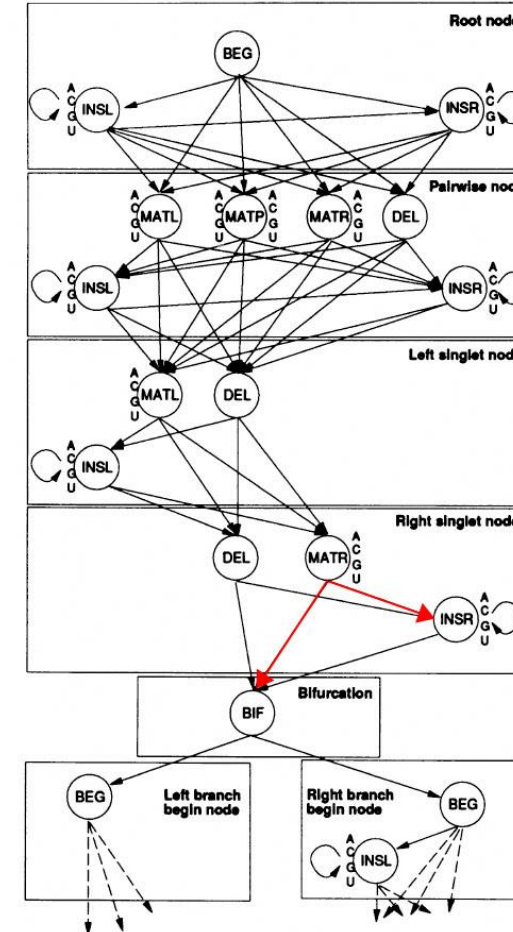


$$S_{i,j,y}(y = MATL, INSL) = \max_{y_{next}} [S_{i+1,j,y_{next}} + \log T(y_{next} | y) + \log \mathcal{P}(x_i | y)]$$

Images—Eddy et al.

Alignment to CM Algorithm

- For each calculation, take into account:
 - Transition (T) to next state.
 - Emission probability (P) in the state as determined by training data.

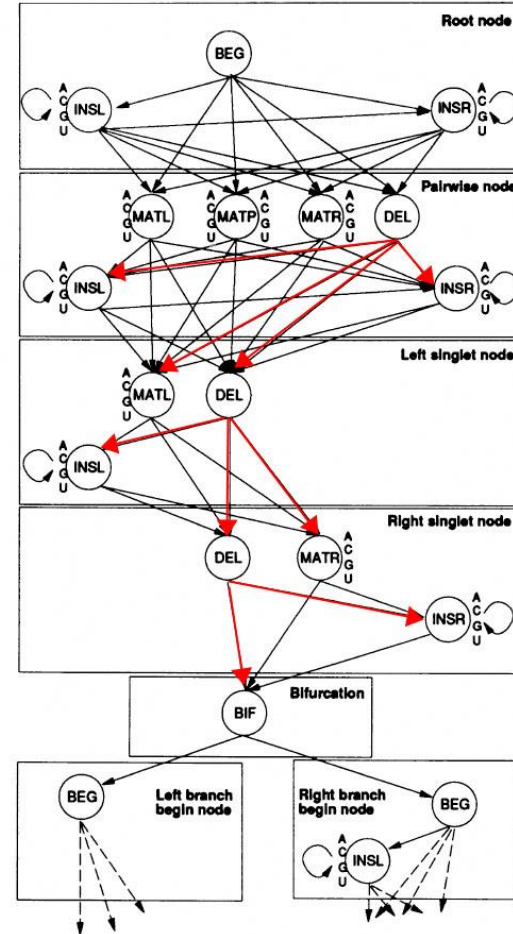


$$S_{i,j,y}(y = MATR, INSR) = \max_{y_{next}} [S_{i,j-1,y_{next}} + \log \mathcal{T}(y_{next} | y) + \log \mathcal{P}(x_j | y)]$$

Images—Eddy et al.

Alignment to CM Algorithm

- For each calculation, take into account:
 - Transition (T) to next state.
 - Emission probability (P) in the state as determined by training data.
- Deletion—does not have emission probability associated with it.

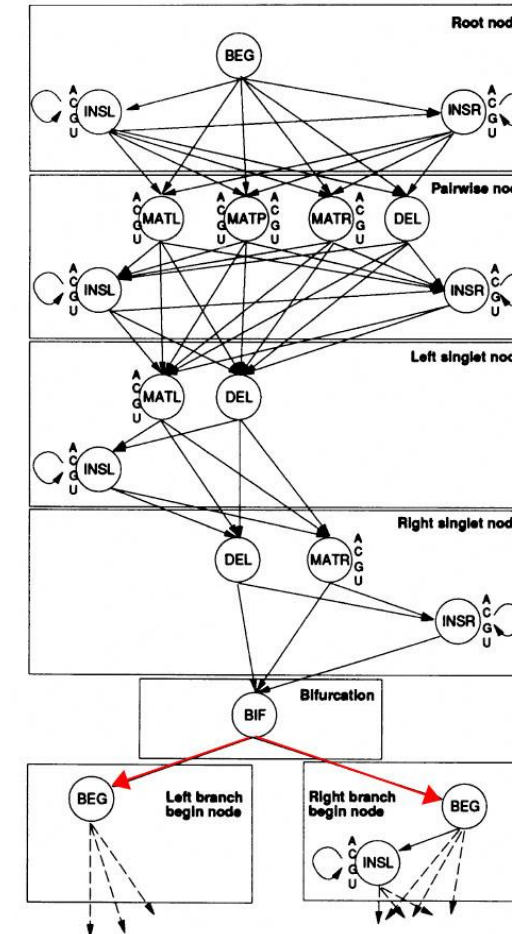


$$S_{i,j,y}(y = DEL) = \max_{y_{next}} [S_{i,j,y_{next}} + \log \mathcal{T}(y_{next} | y)]$$

Images—Eddy et al.

Alignment to CM Algorithm

- For each calculation, take into account:
 - Transition (T) to next state.
 - Emission probability (P) in the state as determined by training data.
- Deletion—does not have emission probability associated with it.
- Bifurcation—does not have state probability associated with it.



$$S_{i,j,y}(y = BIFURC) = \max_{i-1 \leq mid \leq j} [S_{i,mid,y_{left}} + S_{mid+1,j,y_{right}}]$$

Images—Eddy et al.

Covariance Model Drawbacks

- Needs to be well trained.
- Not suitable for searches of large RNA.
 - Structural complexity of large RNA cannot be modeled
 - Runtime
 - Memory requirements