

Combinatorial Optimization in Bioinformatics

RNA Structure Prediction—Part II

Sebastian Will · Yann Ponty

sebastian.will@polytechnique.edu · yann.ponty@lix.polytechnique.fr

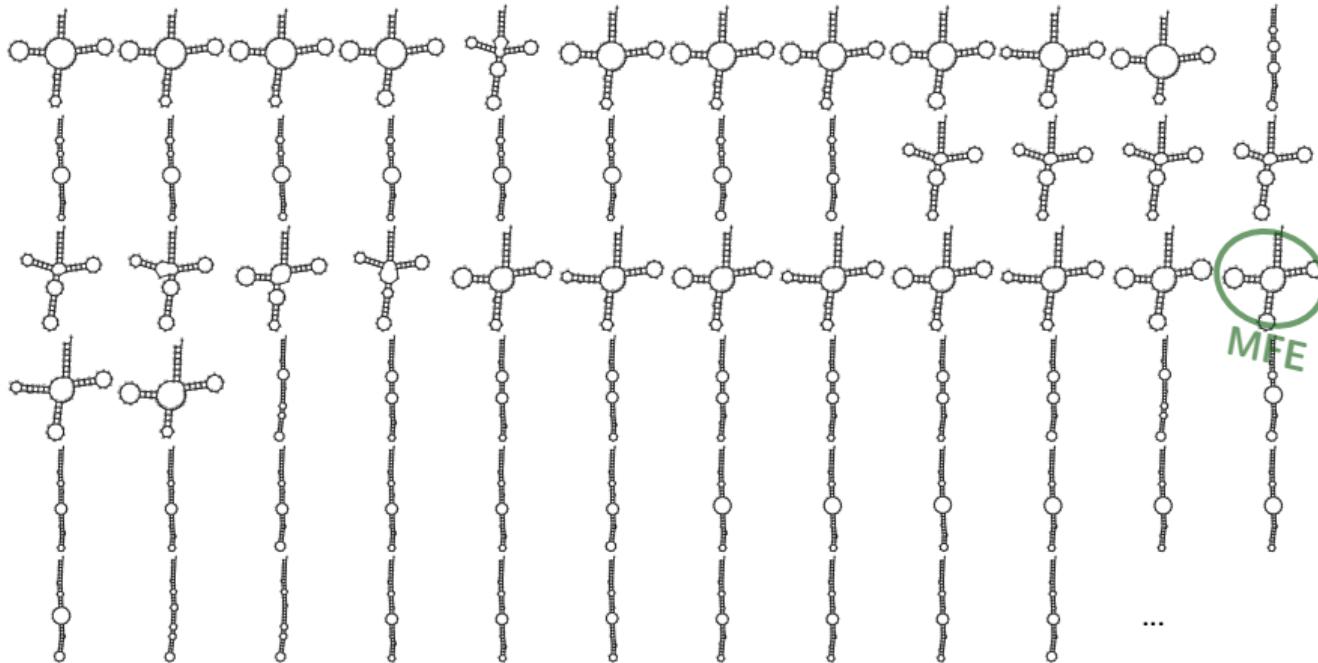


Beyond predicting single RNA structures

- so far: optimization over all (non-crossing) RNA structures
 - from simple to more complex objective functions
 - (further) bioinformatics example of optimization by DP: “optimal substructure”, recursive decomposition
 - first look beyond: relevance of ambiguity, counting
- today: beyond finding the optimal solution—more power of DP
 - **statistical mechanics**: energies correspond to probabilities
 - here: directly **probabilities of structures** in thermodynamic equilibrium, but technique is generally useful to assign probabilities to combinatoric objects.
 - calculate partition functions by DP—**sum over all structures** (in place of optimize)
 - inside + outside algorithm—**probability of features**
 - structure prediction from **multiple sequences, comparative analysis**

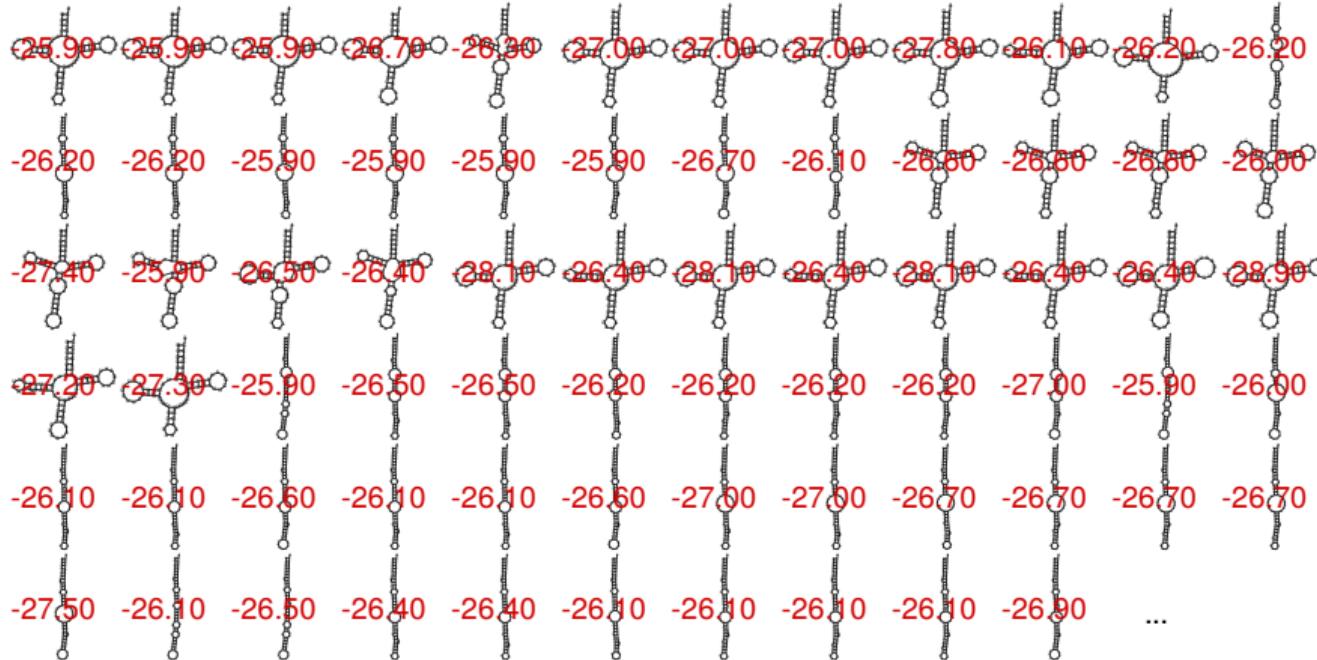
RNA structure ensembles, Probabilities

GGGCUAUUAGCUCAGUUGGUUAGAGCGCACCCUGAUAGGGUGAGGUCGCUGAUUCGAUUCAGCAUAGCCC



RNA structure ensembles, Probabilities

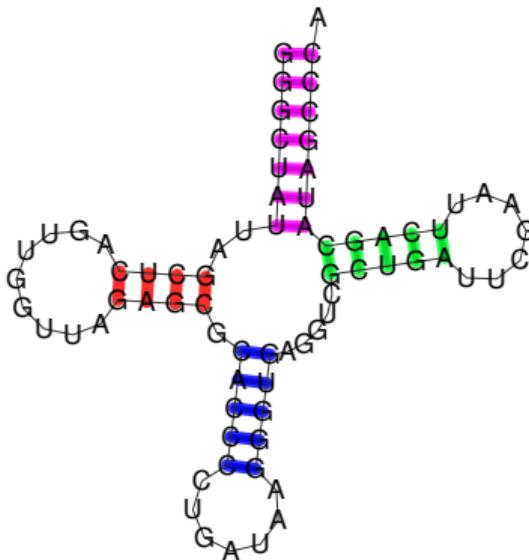
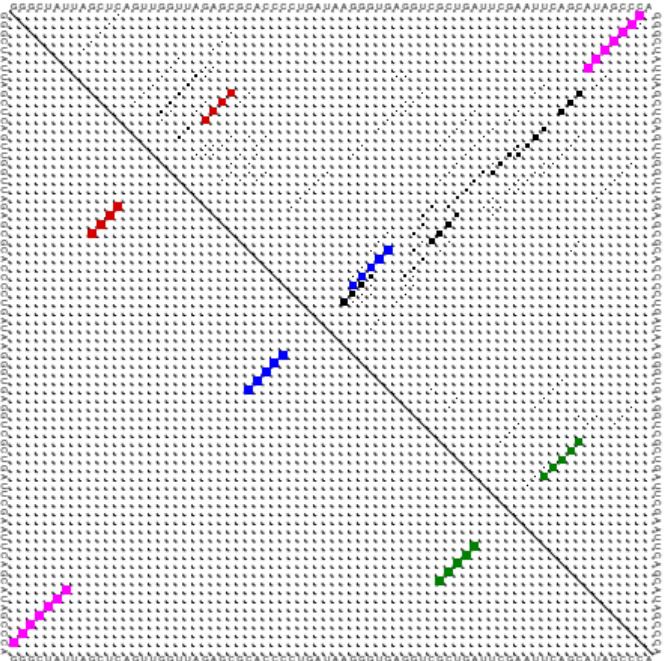
GGGUUUAGCUCAGUUGGUUAGAGCGCACCCUGAUAGGGUGAGGUCGCUGAUUCGAUUCAGCAUAGCCC



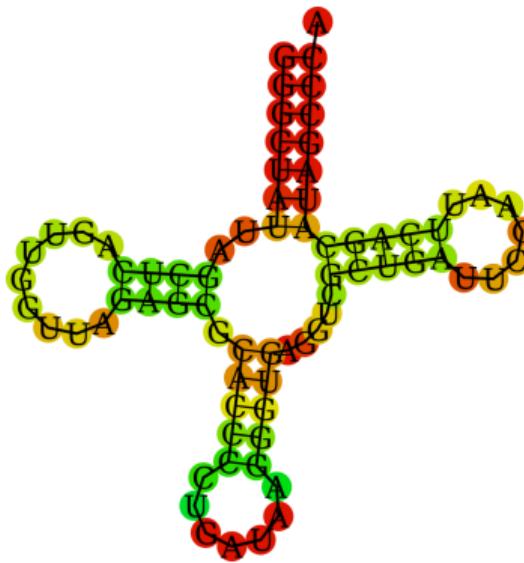
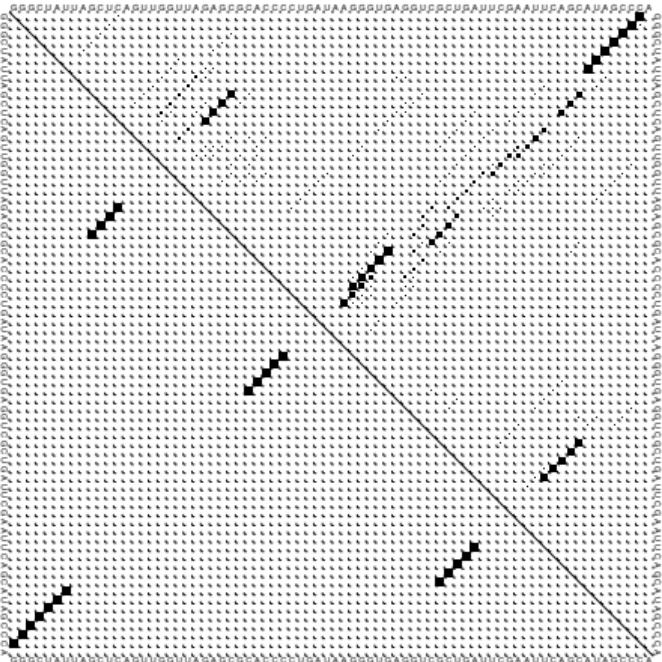
Energies → Structure Probabilities

Structure Probabilities → Base Pair Probabilities

Dotplots and Reliabilities



Dotplots and Reliabilities



Thermodynamic Equilibrium / Statistical Mechanics

We assume that the probability of each possible structure P (in thermodynamic equilibrium) is proportional to its Boltzmann weight (\Rightarrow **Boltzmann distribution**).

$$Pr(P) \propto \exp\left(-\frac{E(P)}{RT}\right)$$

(R = gas constant; T = temperature in Kelvin).

Thermodynamic Equilibrium / Statistical Mechanics

We assume that the probability of each possible structure P (in thermodynamic equilibrium) is proportional to its Boltzmann weight (\Rightarrow **Boltzmann distribution**).

$$Pr(P) \propto \exp\left(-\frac{E(P)}{RT}\right)$$

$(R = \text{gas constant}; T = \text{temperature in Kelvin}).$

Reason to assume **Boltzmann distribution**:

- Formally justified as distribution with maximal entropy in a closed physical system (with constant average energy).
- Boltzmann distributions models the distribution of free, dynamically folding structures after infinite time (and without additional forces).

Normalization: Partition function

$$Pr(P) \propto \exp\left(-\frac{E(P)}{RT}\right)$$

Normalize Boltzmann weights by Z :

$$Pr(P) = \exp\left(-\frac{E(P)}{RT}\right)/Z, \text{ such that}$$

$$\sum_P Pr(P) = 1$$

Normalization: Partition function

$$Pr(P) \propto \exp\left(-\frac{E(P)}{RT}\right)$$

Normalize Boltzmann weights by Z :

$$Pr(P) = \exp\left(-\frac{E(P)}{RT}\right)/Z, \text{ such that}$$

$$\sum_P Pr(P) = 1$$

$$\implies Z = \sum_P \exp\left(-\frac{E(P)}{RT}\right)$$

Z is called partition function.

Normalization: Partition function

$$Pr(P) \propto \exp\left(-\frac{E(P)}{RT}\right)$$

Normalize Boltzmann weights by Z :

$$Pr(P) = \exp\left(-\frac{E(P)}{RT}\right)/Z, \text{ such that}$$

$$\sum_P Pr(P) = 1$$

$$\implies Z = \sum_P \exp\left(-\frac{E(P)}{RT}\right)$$

Z is called partition function.

Efficient computation?

- weighted sum has similarities to counting (DP!)
- weights depend on free energy (Zuker-like DP?!)

Recall: Counting of Structures

$$N_{i,j} = \max \left\{ \begin{array}{l} N_{i+1,j} \\ \max_{\substack{i < k \leq j \\ S[i], S[j] \text{ compl.}}} 1 + N_{i+1,k-1} + N_{k+1,j} \end{array} \right.$$

Recall: Counting of Structures

$$N_{i,j} = \max \left\{ \begin{array}{l} N_{i+1,j} \\ \max_{\substack{i < k \leq j \\ S[i], S[j] \text{ compl.}}} 1 + N_{i+1,k-1} + N_{k+1,j} \end{array} \right.$$

↓

$$\begin{aligned} C_{i,j} = & C_{i+1,j} \\ & + \sum_{\substack{i < k \leq j \\ S[i], S[j] \text{ compl.}}} 1 \cdot C_{i+1,k-1} \cdot C_{k+1,j} \end{aligned}$$

Exchange of operators: $\max \rightarrow \sum / +, \sum / + \rightarrow \prod / \cdot$

From counting to partition function

$$C_{i,j} = C_{i+1,j} + \sum_{\substack{i < k \leq j \\ S[i], S[j] \text{ compl.}}} 1 \cdot C_{i+1,k-1} \cdot C_{k+1,j}$$

$$\Rightarrow Z_{i,j} = Z_{i+1,j} + \sum_{\substack{i < k \leq j \\ S[i], S[j] \text{ compl.}}} \exp\left(\frac{-E_{bp}(i, k)}{RT}\right) \cdot Z_{i+1,k-1} \cdot Z_{k+1,j}$$

From counting to partition function

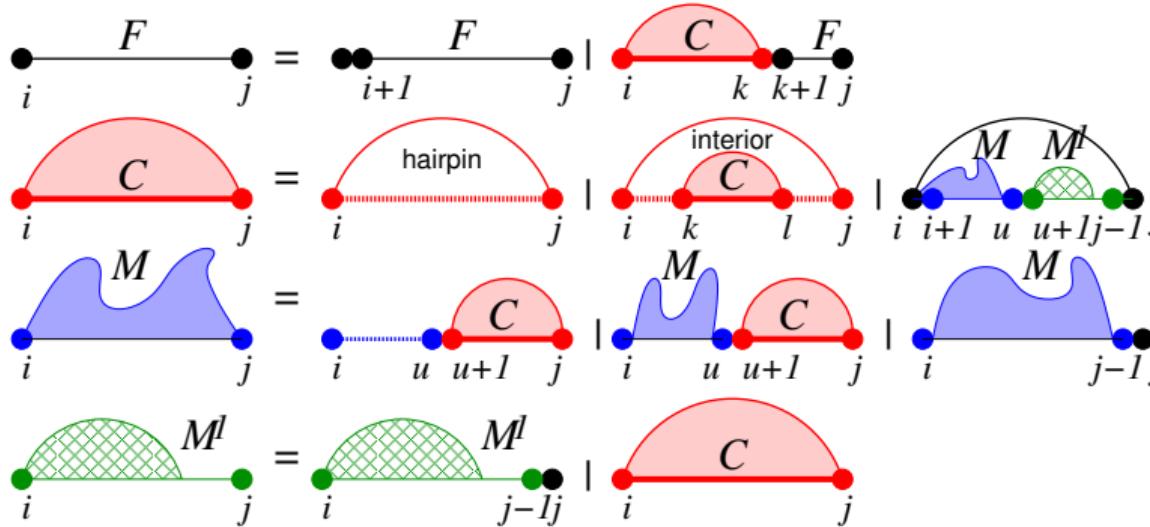
$$C_{i,j} = C_{i+1,j} + \sum_{\substack{i < k \leq j \\ S[i], S[j] \text{ compl.}}} 1 \cdot C_{i+1,k-1} \cdot C_{k+1,j}$$

$$\implies Z_{i,j} = Z_{i+1,j} + \sum_{\substack{i < k \leq j \\ S[i], S[j] \text{ compl.}}} \exp\left(\frac{-E_{bp}(i, k)}{RT}\right) \cdot Z_{i+1,k-1} \cdot Z_{k+1,j}$$

Validity?

$$\begin{aligned} \exp\left(\frac{-E_{bp}(i, k)}{RT}\right) \cdot Z_{i+1,k-1} \cdot Z_{k+1,j} &= \exp\left(\frac{-E_{bp}(i, k)}{RT}\right) \cdot \sum_x \exp\left(\frac{-E(x)}{RT}\right) \cdot \sum_y \exp\left(\frac{-E(y)}{RT}\right) \\ &= \sum_{x,y} \exp\left(\frac{-E_{bp}(i, k)}{RT}\right) \cdot \exp\left(\frac{-E(x)}{RT}\right) \cdot \exp\left(\frac{-E(y)}{RT}\right) \\ &= \sum_{x,y} \exp\left(\frac{-(E_{bp}(i, k) + E(x) + E(y))}{RT}\right) \end{aligned}$$

Free energy minimization ('Zuker' algorithm)



Initialization: $F_{ii} = 0; C_{ii} = M_{ii} = M_{ii}^1 = \infty$

$$\text{Recursions: } F_{ij} = \min \left\{ F_{i+1,j}, \min_{i < k \leq j} C_{ik} + F_{k+1,j} \right\}$$

$$C_{ij} = \min \left\{ H(i,j), \min_{i < k < l < j} C_{kl} + I(i,j; k, l), \min_{i < u < j} M_{i+1,u} + M_{u+1,j-1}^1 + a \right\}$$

$$M_{ij} = \min \left\{ \min_{i < u < j} (u - i + 1)c + C_{u+1,j} + b, \min_{i < u < j} M_{i,u} + C_{u+1,j} + b, M_{i,j-1} + c \right\}$$

$$M_{ij}^1 = \min \{ M_{i,j-1} + c, C_{ij} + b \}$$

Efficient computation of partition functions

Efficient computation: partition function algorithm (McCaskill, 1990) **can be inferred from the MFE algorithm** by *replacing minimum operations with sums and additions with multiplications.*

$$F_{ij} = \min \left\{ F_{i+1,j}, \min_{i < k \leq j} C_{ik} + F_{k+1,j} \right\}$$

$$C_{ij} = \min \left\{ H(i,j), \min_{i < k < l < j} C_{kl} + I(i,j; k, l), \min_{i < u < j} M_{i+1,u} + M_{u+1,j-1}^1 + a \right\}$$

$$M_{ij} = \min \left\{ \min_{i < u < j} (u - i + 1)c + C_{u+1,j} + b, \min_{i < u < j} M_{i,u} + C_{u+1,j} + b, M_{i,j-1} + c \right\}$$

$$M_{ij}^1 = \min \{ M_{i,j-1} + c, C_{ij} + b \}$$

Efficient computation of partition functions

Efficient computation: partition function algorithm (McCaskill, 1990) **can be inferred from the MFE algorithm** by *replacing minimum operations with sums and additions with multiplications.*

$$Z_{i,j} = Z_{i+1,j} + \sum_{i < k \leq j} Z_{i,k}^C Z_{k+1,j}$$

$$Z_{i,j}^C = e^{-\beta H(i,j)} + \sum_{i < k < l < j} Z_{k,l}^C e^{-\beta I(i,j,k,l)} + \sum_{i < u < j} Z_{i+1,u}^M Z_{u+1,j-1}^{M^1} e^{-\beta a}$$

$$Z_{i,j}^M = \sum_{i < u < j} e^{-\beta(u-i+1)c} Z_{u+1,j}^M + \sum_{i < u < j} Z_{i,u}^M Z_{u+1,j}^C e^{-\beta b} + Z_{i,j-1}^M e^{-\beta c}$$

$$Z_{i,j}^{M^1} = Z_{i,j-1}^{M^1} e^{-\beta c} + Z_{i,j}^C e^{-\beta b}$$

(Un)-Ambiguity

Partition function by operator exchange “trick” would not work on ambiguous recursions.
Why?

(Un)-Ambiguity

Partition function by operator exchange “trick” would not work on ambiguous recursions.
Why?

prerequisite, violated by ambiguous decomposition

traceback : structure = 1 : 1

(Un)-Ambiguity

Partition function by operator exchange “trick” would not work on ambiguous recursions.
Why?

prerequisite, violated by ambiguous decomposition

traceback : structure = 1 : 1

Challenge of ambiguity in (original) Zuker-recursions: split of multiloops

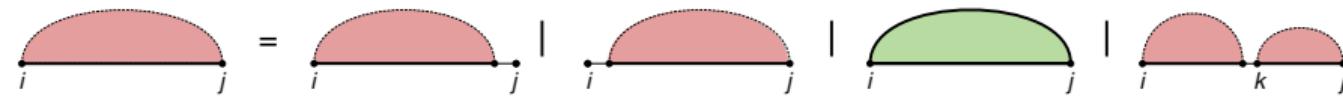
Decomposition of Structure Space by (original) Zuker

$$\text{Diagram 1: } \text{A green semi-circle with endpoints } i \text{ and } j \text{ is equal to the sum of two terms. The first term is a green semi-circle with endpoints } i \text{ and } j \text{ plus a green semi-circle with endpoints } i \text{ and } k \text{ plus a green semi-circle with endpoints } k \text{ and } j.$$

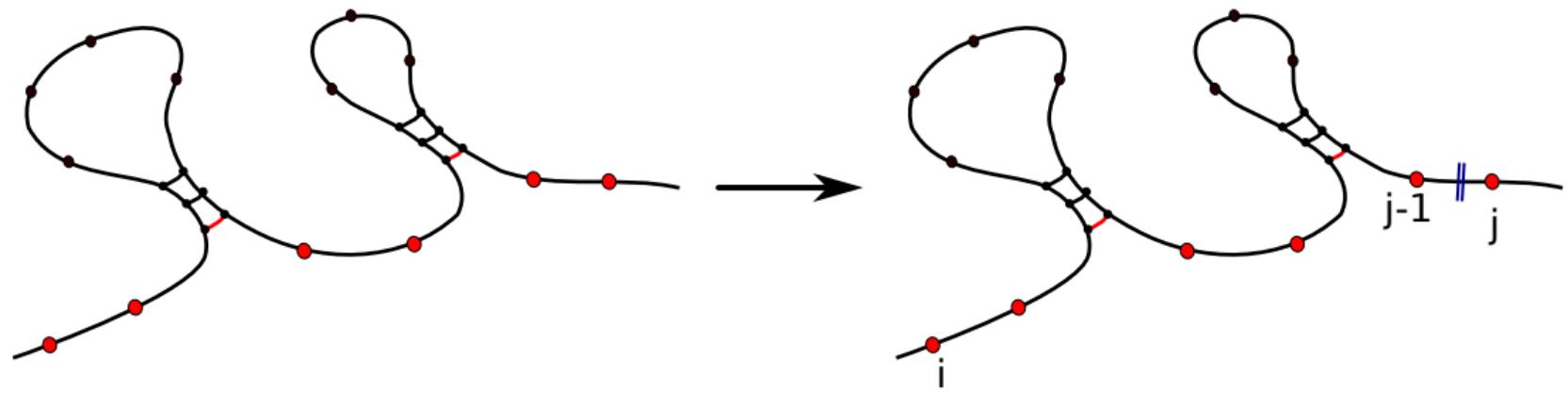
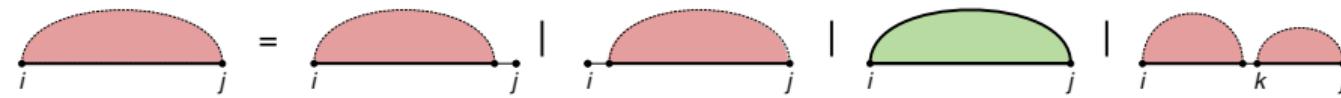
$$\text{Diagram 2: } \text{A green semi-circle with endpoints } i \text{ and } j \text{ is equal to the sum of three terms. The first term is a black semi-circle with endpoints } i \text{ and } j \text{ plus a green semi-circle with endpoints } i \text{ and } i' \text{ plus a green semi-circle with endpoints } j' \text{ and } j.$$

$$\text{Diagram 3: } \text{A red semi-circle with endpoints } i \text{ and } j \text{ is equal to the sum of four terms. The first term is a black semi-circle with endpoints } i \text{ and } j \text{ plus a red semi-circle with endpoints } i \text{ and } i' \text{ plus a green semi-circle with endpoints } i \text{ and } j \text{ plus a red semi-circle with endpoints } i \text{ and } k \text{ plus a red semi-circle with endpoints } k \text{ and } j.$$

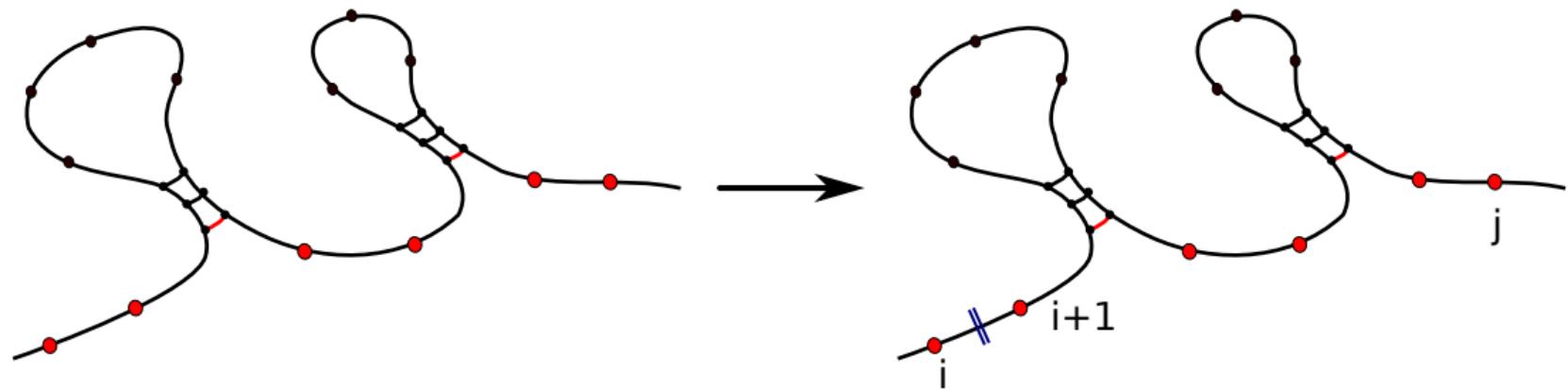
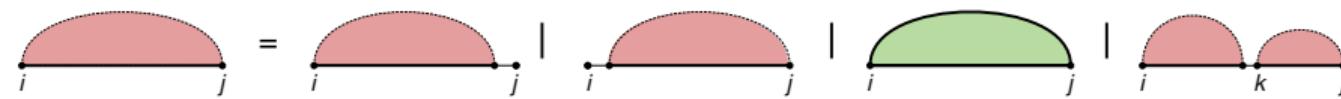
Ambiguity Example



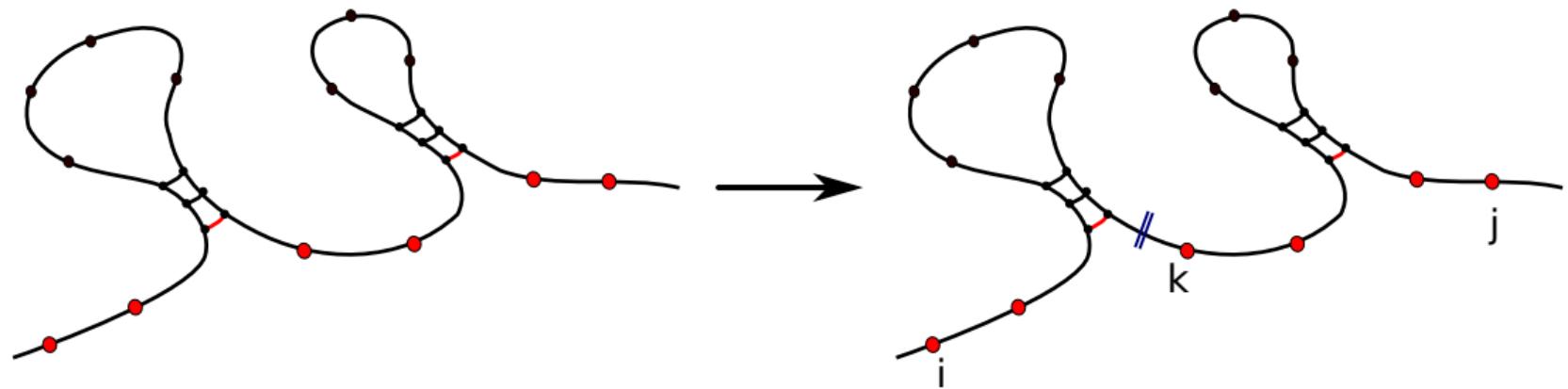
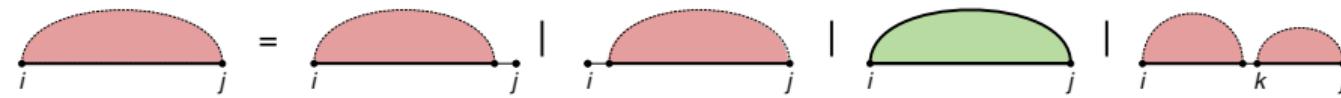
Ambiguity Example



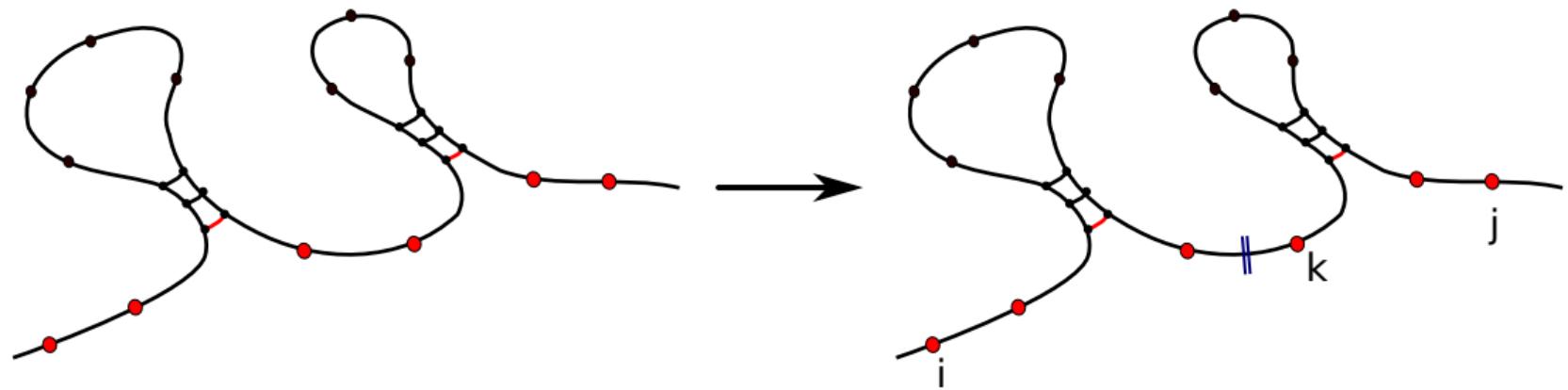
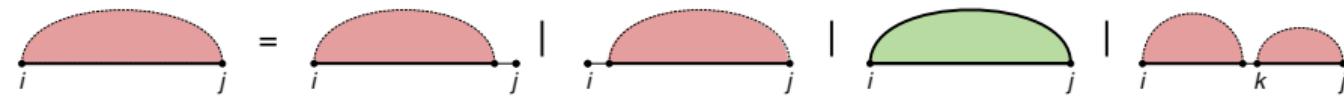
Ambiguity Example



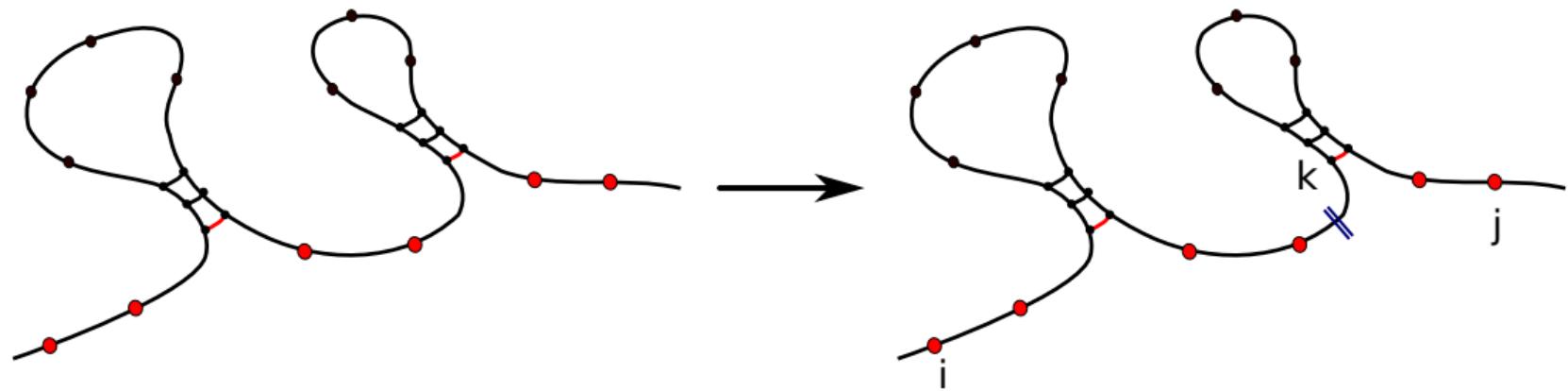
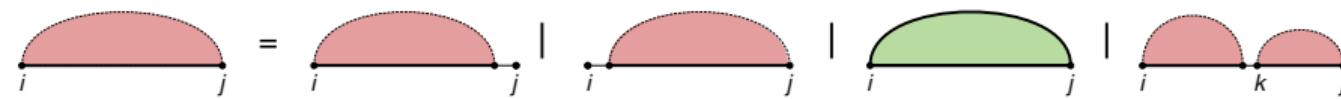
Ambiguity Example



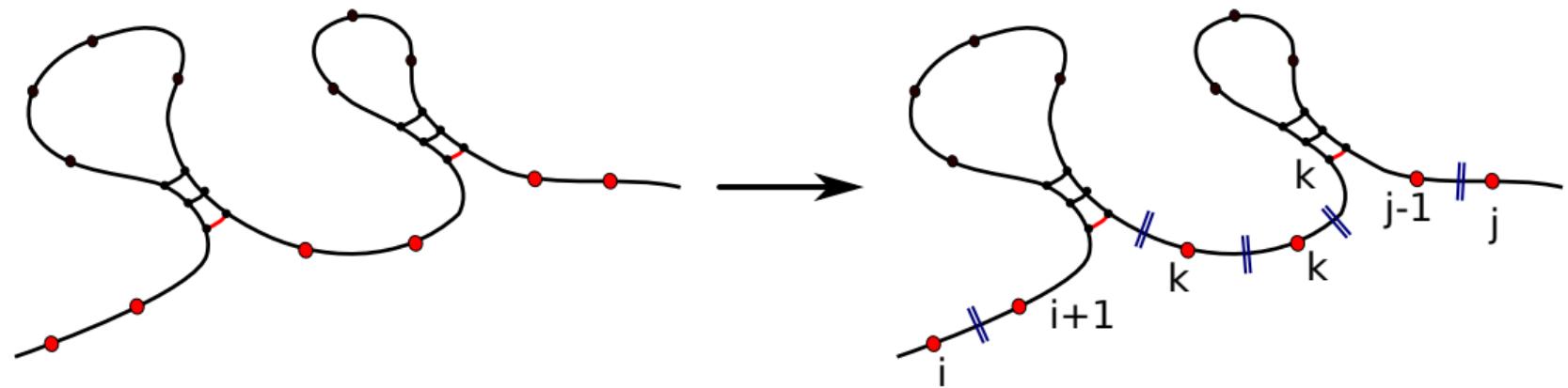
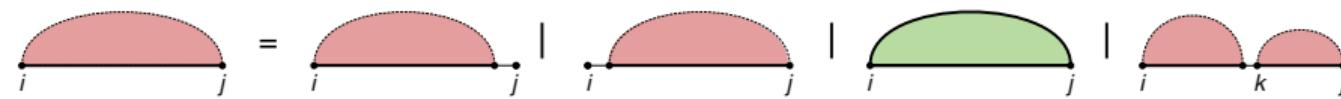
Ambiguity Example



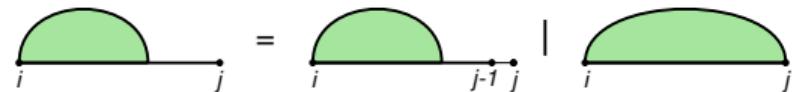
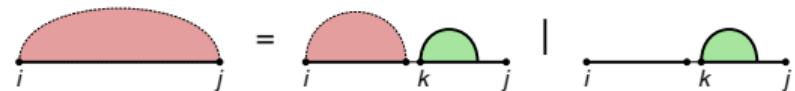
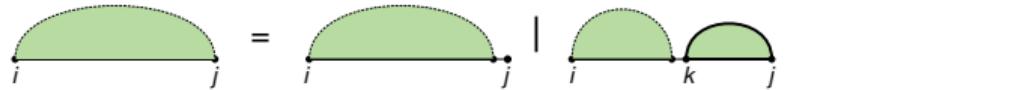
Ambiguity Example



Ambiguity Example



Discarding Ambiguity



Disambiguation due to “fourth” matrix.

Structure and base pair probabilities

The **probability of a specific structure** is given by:

$$Pr(P) = \frac{\exp(-\frac{E(P)}{RT})}{Z}$$

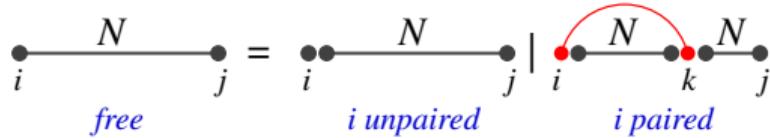
The **probability of a base pair** (i, j) is the fraction of all structures having this base pair and the number of all possible structures.

$$p_{ij} = \sum_{P \ni (i,j)} Pr(P)$$

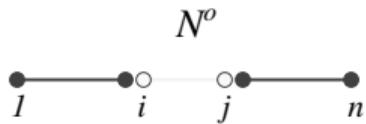
This can be computed efficiently from the product of the **inside** and **outside** partition functions of a base pair:

$$p_{ij} = \frac{Z_{\text{inside}(i,j)} \cdot Z_{\text{outside}(i,j)}}{Z}$$

Outside algorithms

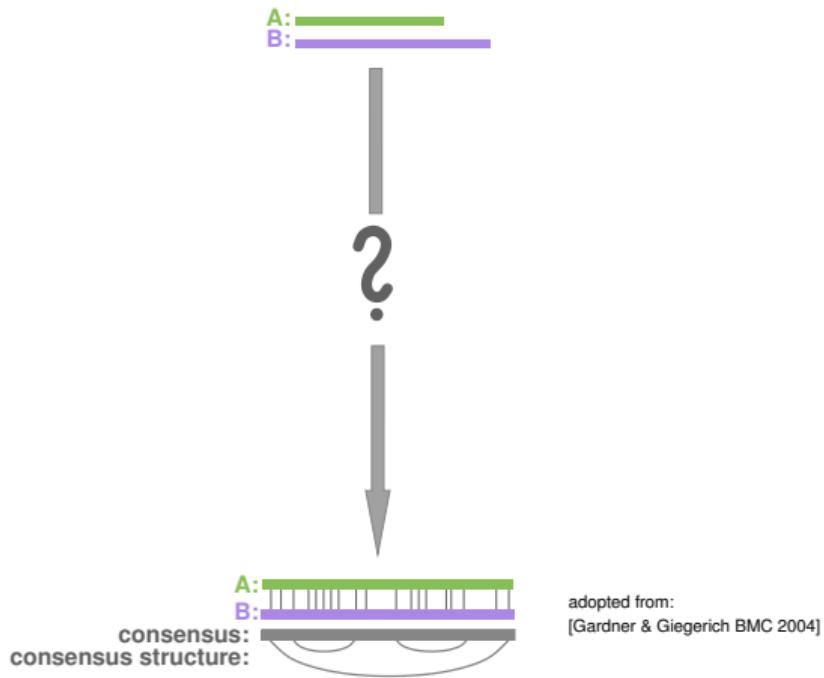


$$N_{i,j} = \max \begin{cases} N_{i-1,j} \\ \max_{1 \leq k < j} 1 + N_{i+1,k-1} + 1 + N_{k+1,j} : S_i \text{ and } S_k \text{ compl.} \end{cases}$$



$$N_{i,j}^o = \max \begin{cases} N_{i-1,j}^o & i - 1 \text{ unpaired} \\ \max_{\substack{1 \leq k < i \\ S_k, S_{i-1} \text{ compl.}}} N_{k,j}^o + 1 + N_{k+1,i-2} & i - 1 \text{ paired left} \\ \max_{\substack{j \leq k < n \\ S_{i-1}, S_k \text{ compl.}}} N_{i-1,k}^o + 1 + N_{j,k-1} & i - 1 \text{ paired right} \end{cases}$$

Comparative RNA Analysis



adopted from:
[Gardner & Giegerich BMC 2004]

Typical Scenario

Given: set of related RNA sequences

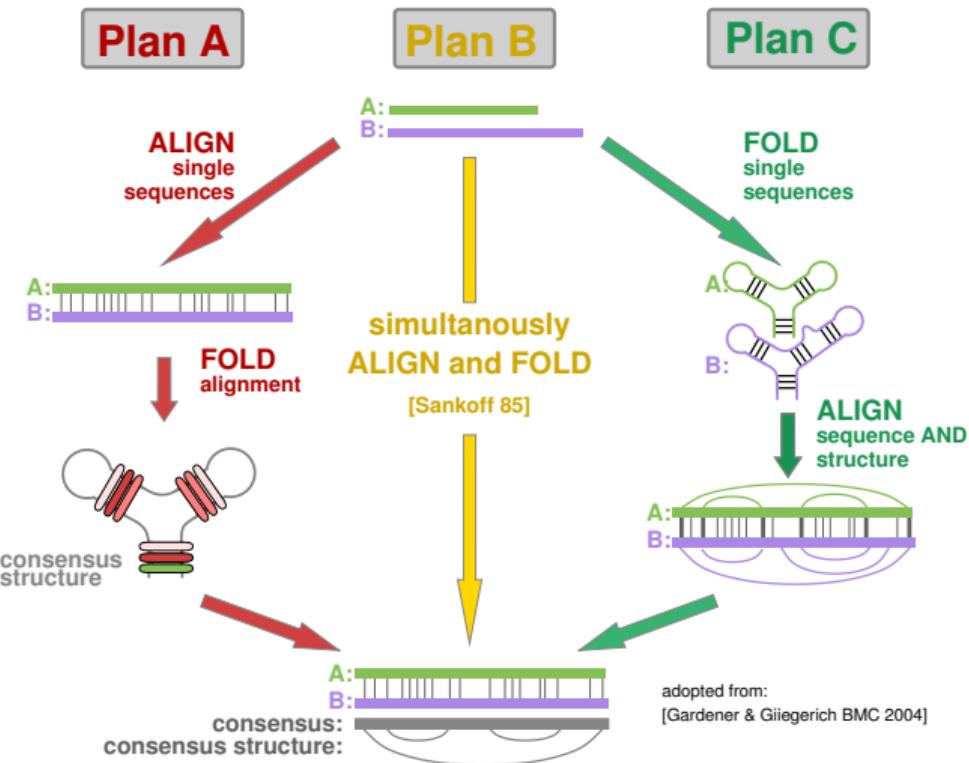
```
>AF008220  
GGAGGAUUAGCUCAGCUGGGAGAGCAUCUGCCUACAAGCAGAGGGUCGGCGGUUCGAGCCCGUCAUCCUCA  
>M68929  
GCGGAUAUAACUUAGGGUAAAAGUUGCAGAUUGUGGCUCUGAAACACGGGUUCGAAUCCGUUAUUCGCC  
>X02172  
GCCUUUAUAGCUUAGUGGUAAAGCGAUAACUGAAGAUUUUUACAUGUAGUUCGAUUCUCAUUAAGGGCA  
>Z11880  
GCCUUCUAGCUCAGUGGUAGAGCGCACGGCUUUUAACCGUGUGGUUCGGAUCCCCACCGGAAGGCG  
>D10744  
GGAAAAUUGAUCAUCGGCAAGAUAGUUUUACUAAAUAAGGAUUUAACCUUGGUGAGUUCGAAUCUCACAUUUUCG
```

Wanted:

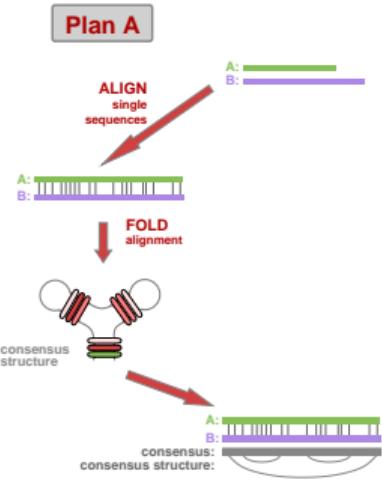
```
AF008220 GGAGGAUU-AGCUCAGCUGGGAGAGCAUCUGCCUACAAGC-----AGAGGGUCGGCGGUUCGAGCCCGUCAUCCUCA  
M68929 GCGGAUAU-AACUUAGGGUAAAAGUUGCAGAUUGUGGCUC-----UGAAAAA-CACGGGUUCGAAUCCGUUAUUCGCC  
X02172 GCCUUUAU-AGCUUAG-UGGUAAAGCGAUAACUGAAGAUU-----UAUUUACAUGUAGUUCGAUUCUCAUUAAGGGCA  
Z11880 GCCUUCU-AGCUCAG-UGGUAGAGCGCACGGCUUUUAACC-----GUGUGGUUCGUGGUUCGGAUCCCCACCGGAAGGCG  
D10744 GGAAAAUUGAUCAUCGGCAAGAUAGUUUUACUAAAUAAGGAUUUAACCUUGGUGAGUUCGAAUCUCACAUUUUCG  
  
consens (((((((.....))))(((((.....)).....))))....((((.....)))))))))).
```

- Typically, only sequences are known, not their structures
- Evolutionary conserved for many RNAs: sequence **and** structure

RNA Alignment Approaches



Comparative RNA Analysis: Plan A



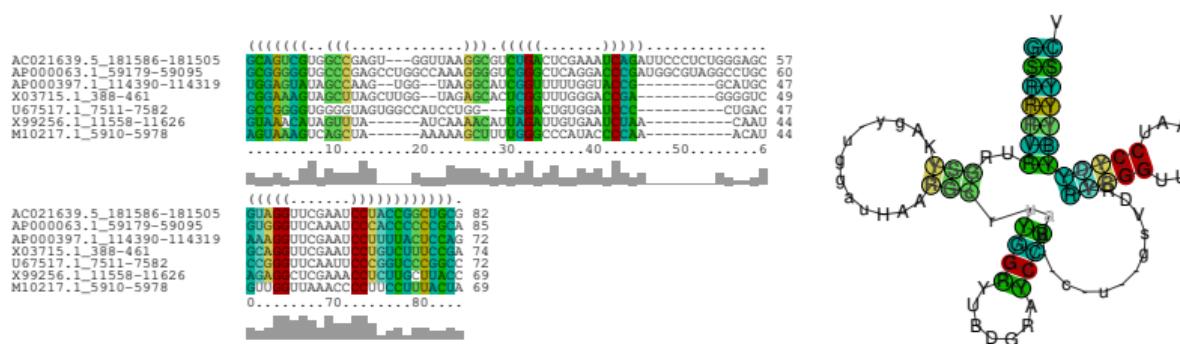
Remarks

- **ALIGN:** multiple sequence alignment
(recall Needleman-Wunsch, progressive alignment)
- **FOLD:** generalize prediction for single sequences

Predict (optimal) consensus structure of an RNA alignment

CLUSTAL W --- LocARNA 2.0.0RC7

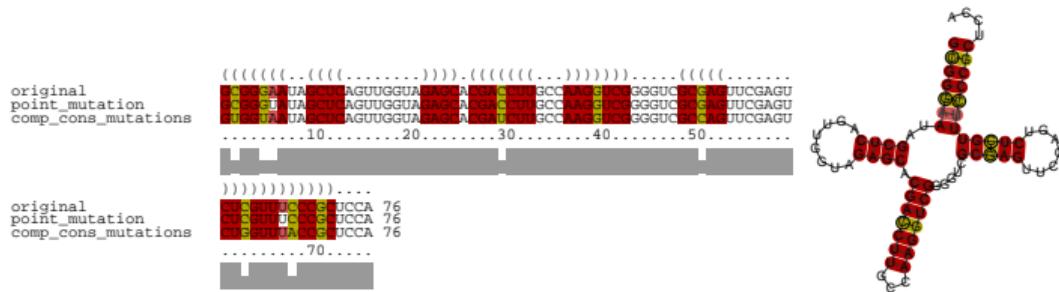
AC021639.5_181586-181505 GCAGUCGUGGCCGAGU---GGUUAAGGCCUGACUCGAAUCAGAUCCCCUUGGGAGCGUAGGUUCGAUCCUACCGGCUGCC
AP000063.1_59179-59095 GCGGGGGUGGCCGAGCCUGGCCAAGGGGUCGGCUCAGGACCCGAUGGCGUAGGCCUGCGUGGGUUCAAUCCCACCCCCCGCA
AP000397.1_114390-114319 UGGAGUAUAGCCAAG--UGG--UAAGGCAUCGGUUUUUGGUACG-----GCAUGCAAAGGUUCGAUCCUUUUACUCCAG
X03715.1_388-461 CGGAAAGUAGCUUAGCUUUGG--UAGAGCACUCGGUUUGGGACCGA-----GGGGUCGCCAGGUUCGAUCCUGUCUUCCGA
U67517.1_7511-7582 GCGGGGGUGGGGUAGUGGCCAUCCUGG---GGGACUGUGGAUCCC-----CUGACCCCCGUCAUCCCCGGUCCGGCC
X99256.1_11558-11626 GUAAAACAUAGUUUA-----AUCAAAACAUAGAUUUGUGAAUCUAA-----CAAUAGAGGCCUCAACCUCUUGCUUACCC
M10217.1_5910-5978 AGUAAAAGUCAGCUA-----AAAAAGCUUUUGGGCCAUACCCCAA-----ACAUGUUGGUAAAACCCCUUCCUUACUAA



RNAalifold

IN: RNA alignment

OUT: **Consensus structure of aligned RNAs**

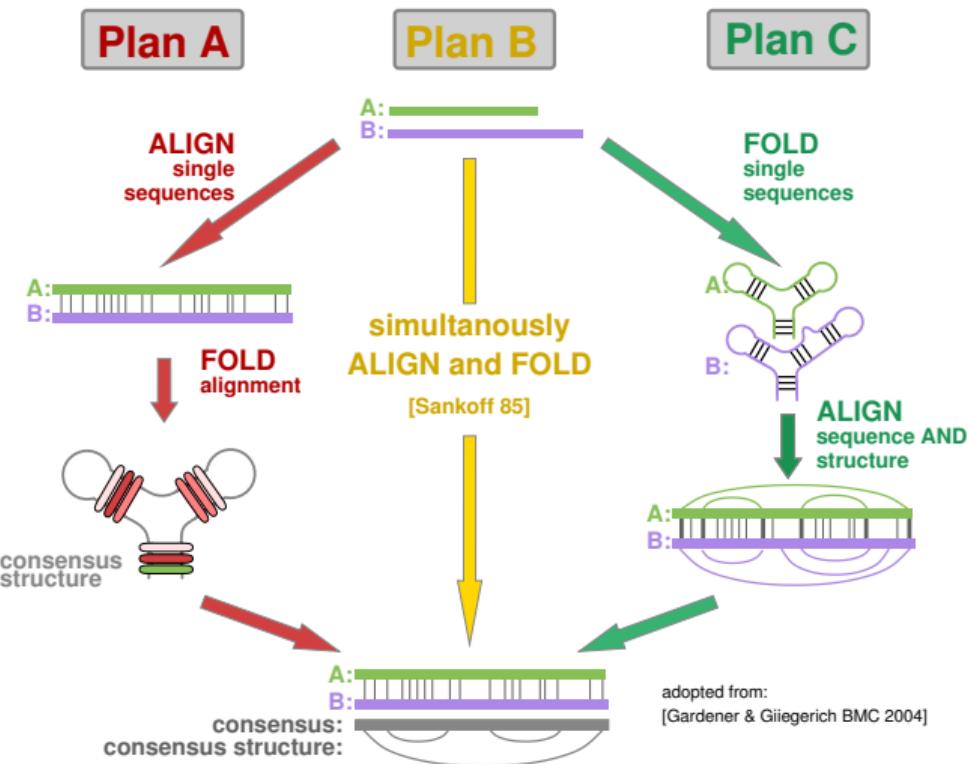


Optimizes **free energy** + **conservation score**

conservation score = **compensatory and consistent mutations**

- alignment as sequence of alignment columns. Folding of this sequence is analogous to folding of an RNA sequence, 'base pairs' are formed between alignment columns
- Thus, same decomposition as Zuker; but modified scoring: sum of loop energies over all sequences & add conservation score

RNA Alignment Approaches

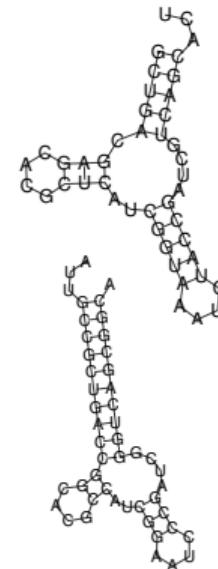
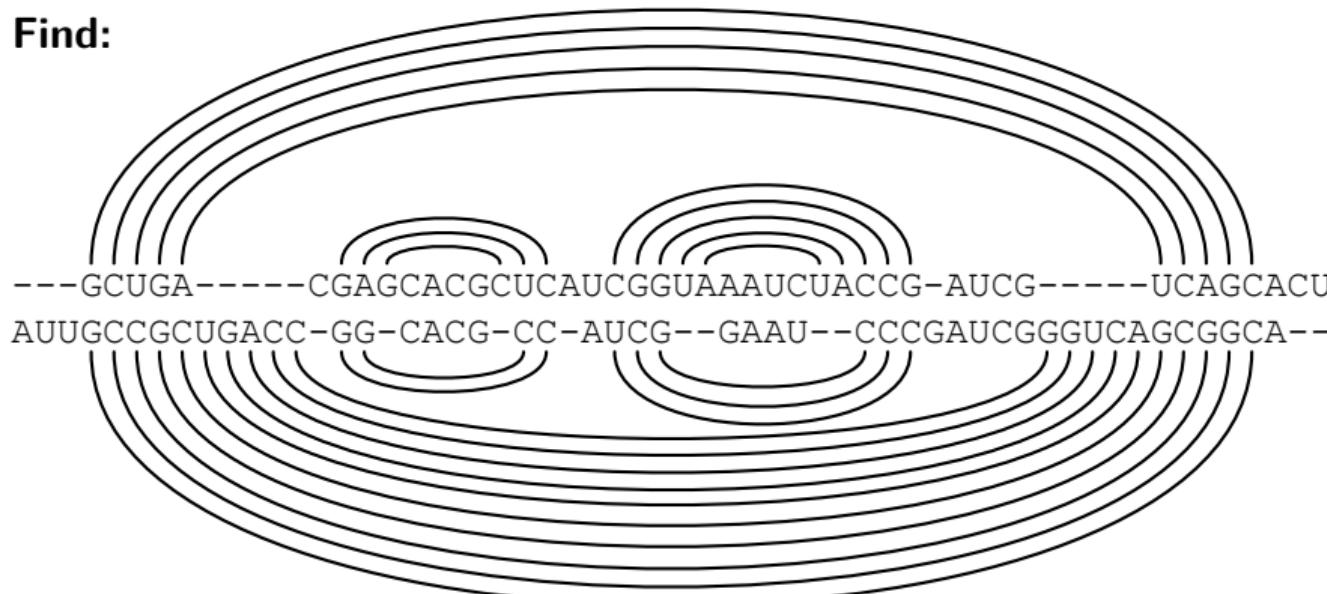


Simultaneous Alignment and Folding [Sankoff]

Given: $A = \text{GCUGACGAGCACGCUCAUCCGUAAAUCUACCGAUCGUCAGCACU}$

& $B = \text{AUUGC CGCUGACC GG CACG CC AUCG -- GAAU -- CCCGAUCGGUCAGCGGCA}$

Find:



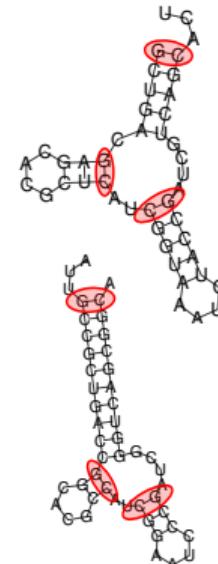
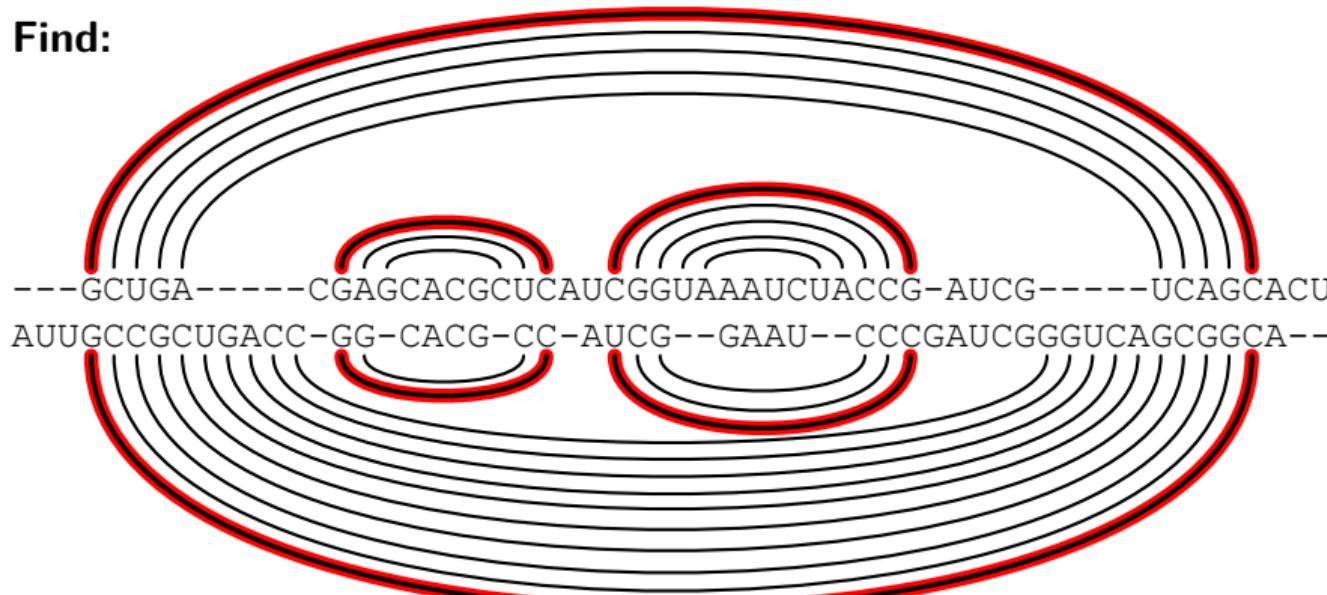
edit distance + energy A + energy B \rightarrow opt

Simultaneous Alignment and Folding [Sankoff]

Given: $A = \text{GCUGACGAGCACGCUCAUCCGUAAAUCUACCGAUCGUCAGCACU}$

& $B = \text{AUUGC CGC UGACC GGG ACGCC AUCGG AAUCC CGAUCGG UCAGCGGCA}$

Find:



edit distance + energy A + energy B \rightarrow opt

Sankoff Problem Definition

- Idea: Sankoff = Zuker Folding \times Needleman/Wunsch Alignment
- IN: two sequences a and b
- find two **equivalent** structures P_a and P_b
and **compatible** alignment A of a and b
such that $Energy(a, P_a) + Energy(b, P_b) + EditDistance(A)$ minimal
- where: $Energy$ yields (loop-based) Turner free energy,
 $EditDistance$ is edit distance (base mismatch x, indel y)
- what means **compatible**
 A must align all base pairs $(i_f, i_g) \in P_a$ and $(j_f, j_g) \in P_b$ that enclose branches
(red base pairs of previous slide)

Constraints

Summarizing, we need to find the optimal structures + alignment with the following constraints:

constraints on the **predicted structures**:

- must be equivalent (same shape)

constraints on the **alignment**:

- multiloops must be aligned to their equivalent partner
- hairpin loops must be aligned to their equivalent partner
- each 2-loop (includes stacking and bulge) must be aligned to exactly one other 2-loop or must be aligned to gaps entirely.

Edit distance of sub-sequences

- distance based score
 - $x = \text{base mismatch}$
 - $y = \text{base deletion}/\text{insertion}$
- $D(i_1, j_1; i_2, j_2)$ minimum sequence alignment cost between sequences $a_{i_1} \dots a_{j_1}$ and $b_{i_2} \dots b_{j_2}$.

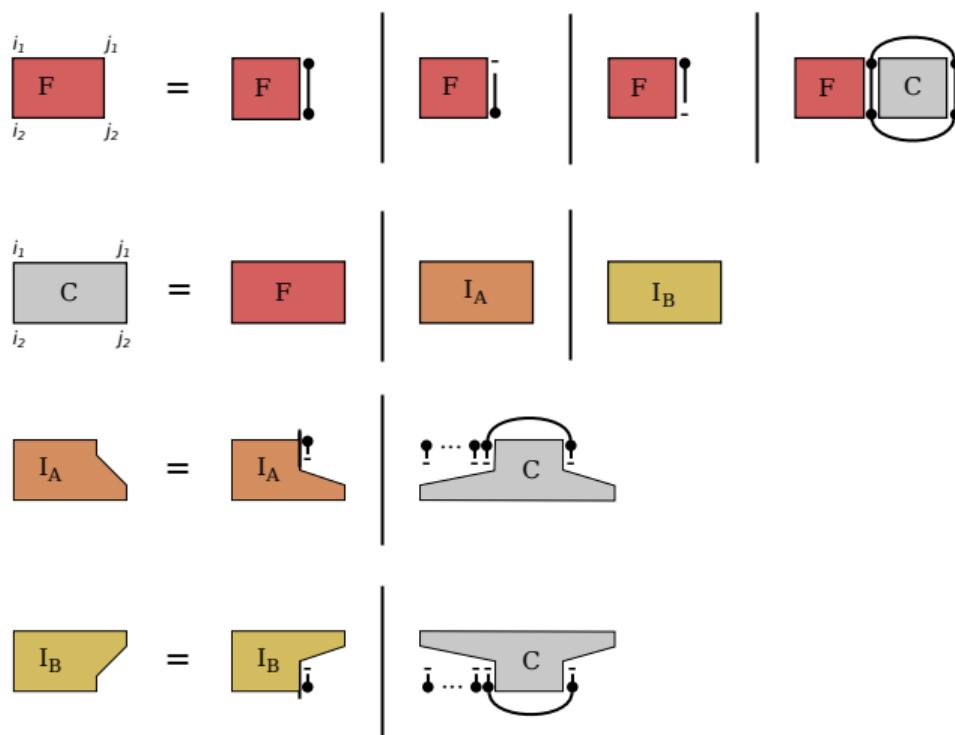
- Recursion: $D(i_1, j_1; i_2, j_2) = \min \begin{cases} D(i_1, j_1 - 1; i_2, j_2 - 1) + \begin{cases} x & \text{if } a_{j_1} \neq b_{j_2} \\ 0 & \text{otherwise} \end{cases} \\ D(i_1, j_1 - 1; i_2, j_2) + y \\ D(i_1, j_1; i_2, j_2 - 1) + y \end{cases}$
- Initialization: $D(i_1, i_1; i_2, i_2) = \begin{cases} x & \text{if } a_{i_1} \neq b_{i_2} \\ 0 & \text{else} \end{cases}$

Simplified Sankoff (without full energy model)

We define the following functions (which will be evaluated recursively):

- $F(i_1, j_1, i_2, j_2)$ best alignment between subsequences $a_{i_1} \dots a_{j_1}$ and $b_{i_2} \dots b_{j_2}$.
- $C(i_1, j_1, i_2, j_2)$ best “closed” alignment between the same subsequences, where we predict the base pairs (i_1, j_1) and (i_2, j_2) .
- $I_A(i_1, j_1, i_2, j_2)$ alignment of the same subsequences, which deletes j_1
- $I_B(i_1, j_1, i_2, j_2)$ symmetrically (...which inserts j_2)

Simplified Sankoff (without full energy model)



Sankoff's extreme complexity

space complexity $O(n^4)$

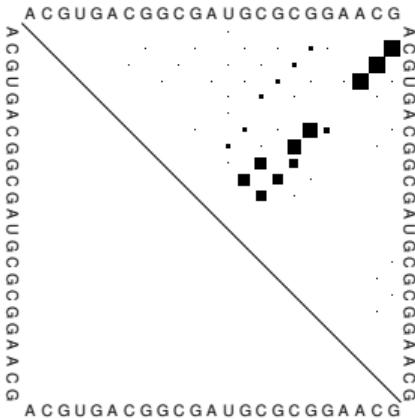
- constant number of matrices
- each of them has $O(n^4)$ entries

time complexity $O(n^6)$

- each entry of matrix D requires constant time
- each entry of F, C, and G requires $O(n^2)$ time (minimize over all h_1, h_2)
- hence: $n^4 \cdot n^2 = n^6$

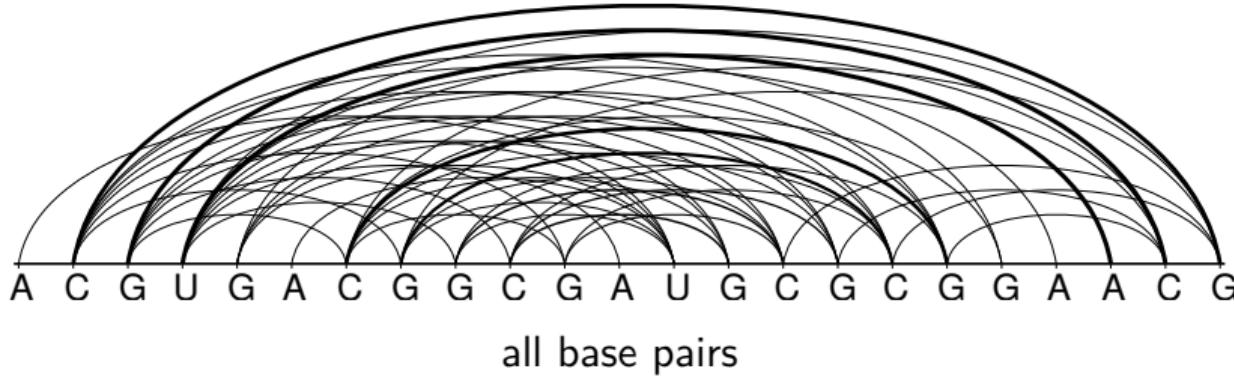
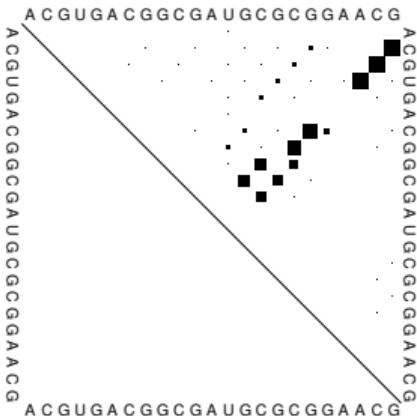
Ensemble-based sparsification (LocARNA)

- Sparsify structure ensemble



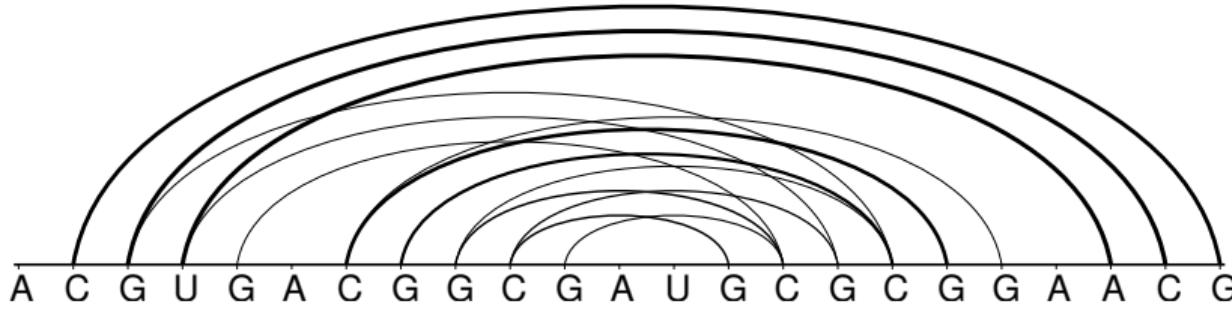
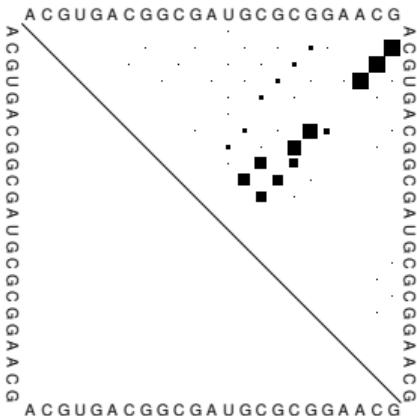
Ensemble-based sparsification (LocARNA)

- Sparsify structure ensemble



Ensemble-based sparsification (LocARNA)

- Sparsify structure ensemble



only probable base pairs

Ensemble-based sparsification (LocARNA)

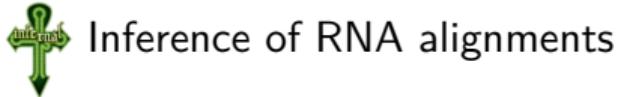
- Sparsify structure ensemble



- **immediate time improvement** by $O(n^2)$
- **space improvement** by $O(n^2)$ requires rearrangement of recursion evaluation

Rfam / Infernal

- Infernal: scan genomic data for RNA family members



- important tool for Rfam

Rfam 10.1 (June 2011, 1973 families)

<http://rfam.sanger.ac.uk/>

- in Rfam: 'hand-curated' seed alignments \Rightarrow full alignments
- use Stochastic Context Free Grammars to model RNA families
- model of a family: *Consensus Model (CM)*

input multiple alignment:

[structure] . : : <<< >- >>: <<- < . >>> .

human . AAGACUUCGGAUCUGGCG . A C A . CCC .

mouse aUACACUUUCGGAUG - CACC . AAA . GUG a

orc . AGGUUCUUC - GCACGGGCAgCCA c UUC .

1 5 10 15 20 25 28

example structure:

U C
U G₁₀
C G A

5 A U

G C

A U¹⁵

27 G G C G A²¹

2 A

• •

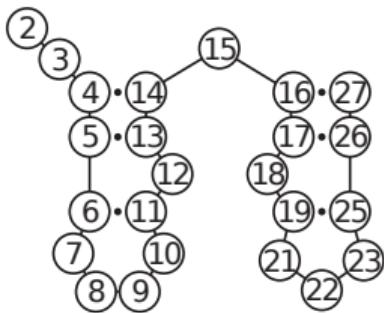
C C

25 C A

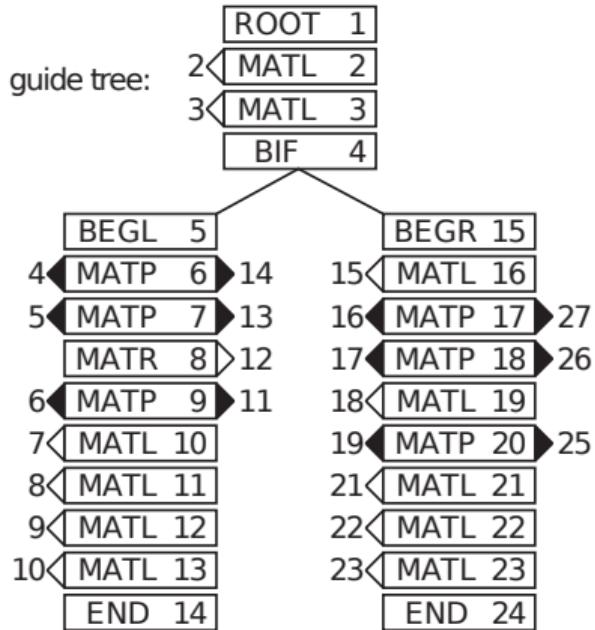
Infernal

Construct grammatical description

consensus structure:



guide tree:



Infernal

- Construct CM from guide tree
- Expand nodes of guide tree:
 - Add match, insertion, and deletion states
- learn transition and output probabilities from alignment
- CM comparable to profile HMM for protein families (Pfam)

