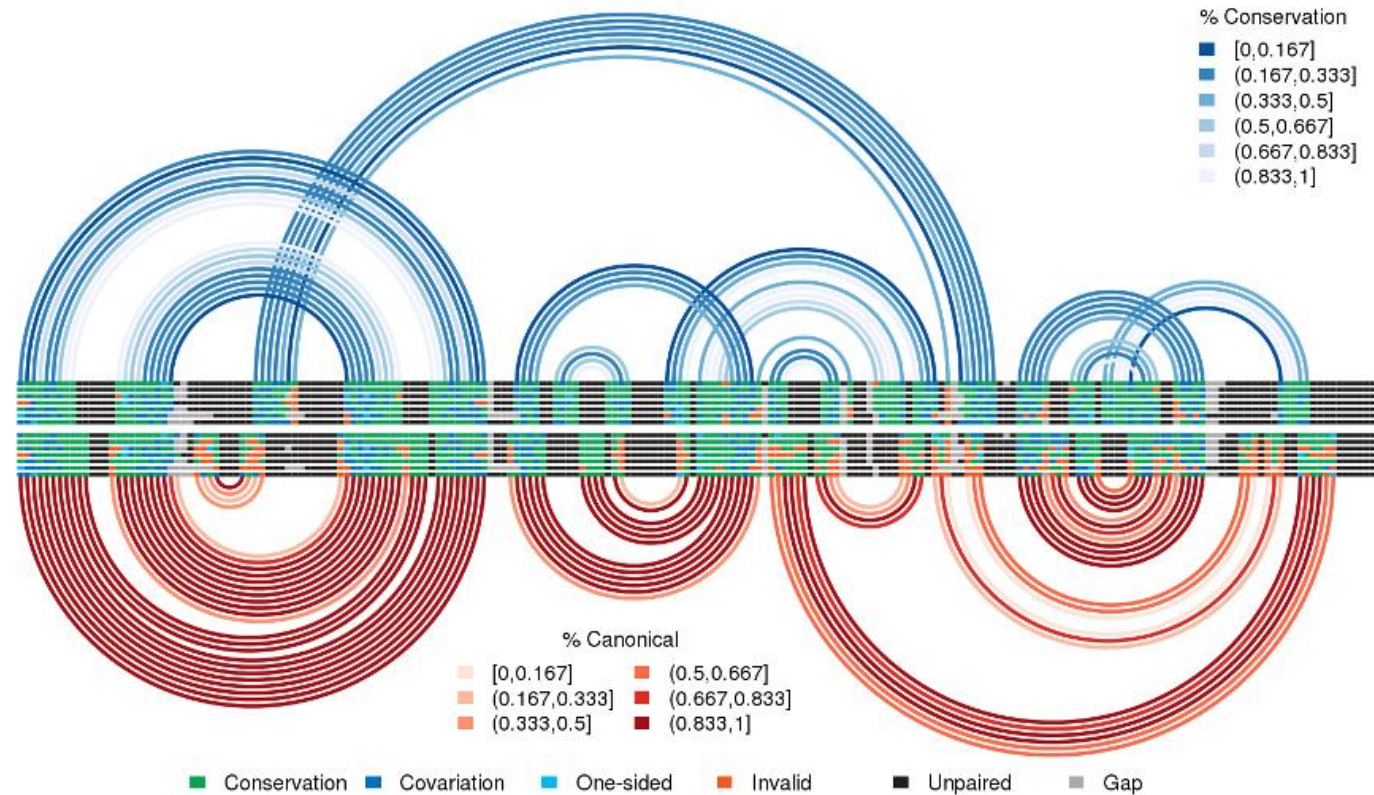


# EECS730: Introduction to Bioinformatics

## Lecture 10: Non-coding RNA secondary structure alignment

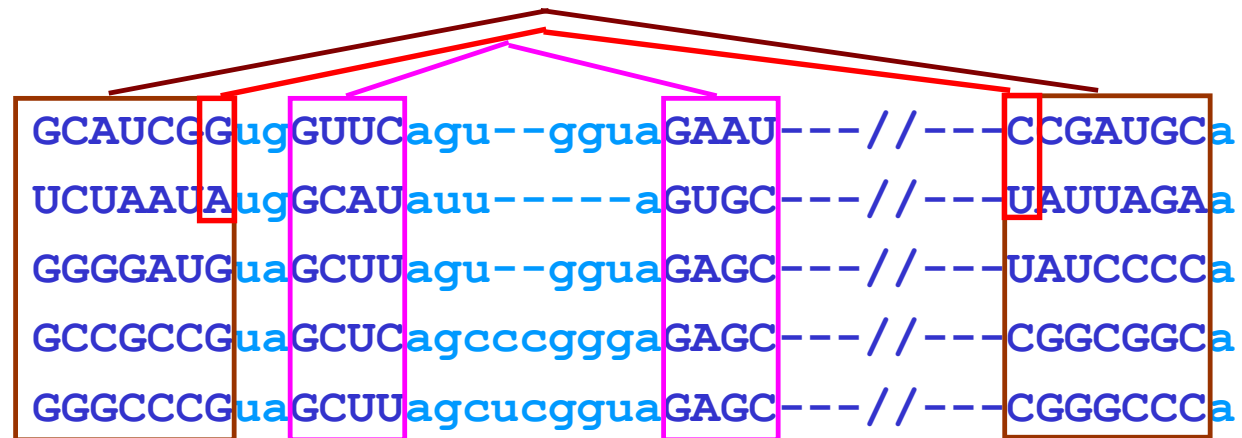
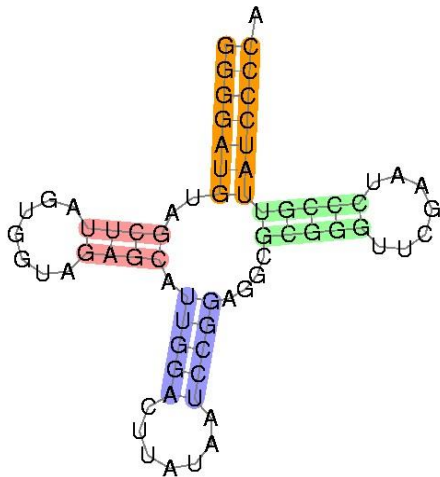


<http://www.e-rna.org/r-chie/images/doublecov.png>

Some slides were adapted from Dr. Shaojie Zhang (University of Central Florida)

# RNA secondary structure alignment

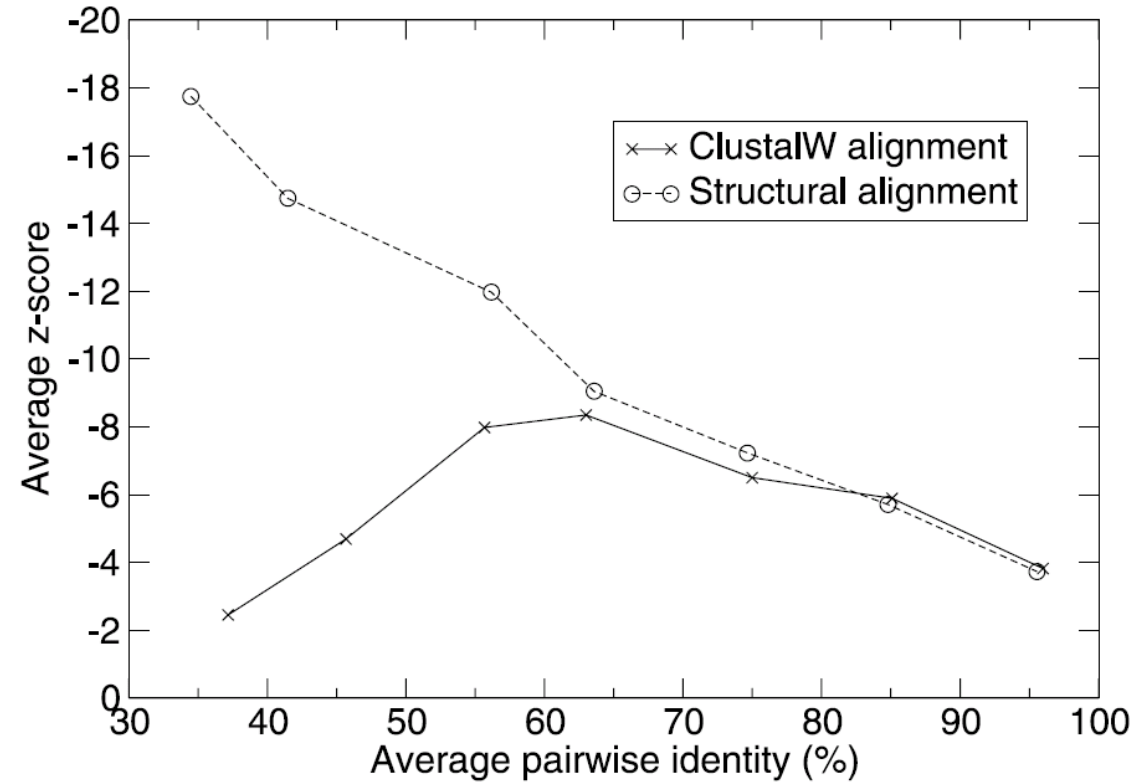
- Sequence alignment will not work for RNAs with low sequence similarity
- RNA conserved on secondary structure rather than primary sequence



# Failure of sequence alignment on RNA

```
ACCCG-UUAAU
|  |||  ||*||
A-CCGUUUCAU
```

```
Stru1: >>>>  <<<<<
Seq1:  GGGGCAACCCC
      +++++* ||+++++
Seq2:  AUCCGAAGGAU
      | | | | |
      | | | | |
      | | | | |
      | | | | |
      | | | | |
```



(Based on pairwise alignments of SRP RNAs)

# Need to consider secondary structure while comparing RNAs

- Sequence-Sequence alignment
  - Only works for high sequence similarity RNAs
- Structure-Structure alignment
  - Compute similarity of two RNA structures
- Sequence-Structure alignment
  - Can the sequence fold into a given secondary structure
- Simultaneous Alignment and folding
  - Compute the consensus structure two RNA sequences can fold into

# Sequence-Sequence alignment

- Will only work for high-sequence similarity RNAs
- Can incorporate secondary structure information into sequence by designing a new alphabet
- “For each piece, RNAfold predicts whether each nucleotide is paired upstream, paired downstream or unpaired. To take advantage of fast primary sequence homology search programs, we map these sequences into a 12 letter alphabet representing nucleotide plus pairing direction.”

# Sequence-Sequence alignment (folded-BLAST)

Table 2. CM scan best recovery motif comparison. For each ncRNA family and each homology search program used, the motif/CM that recovered the most instances of the particular family is listed. The actual motif identifications can be cross-referenced online. *sen.* is the recovery percentage (sensitivity), and *spe.* is the specificity of the CM scan; “None” indicates that no instances were recovered.

ncRNA family	NCBI-BLAST		WU-BLAST		SSEARCH		<i>folded</i> -BLAST	
	<i>sen.</i>	<i>spe.</i>	<i>sen.</i>	<i>spe.</i>	<i>sen.</i>	<i>spe.</i>	<i>sen.</i>	<i>spe.</i>
t-box	0.69	0.98	<b>0.71</b>	<b>0.99</b>	0.41	<b>0.99</b>	0.68	<b>0.99</b>
SAMI	0.94	<b>0.99</b>	<b>0.96</b>	<b>0.99</b>	0.84	<b>0.99</b>	0.94	<b>0.99</b>
TPP	0.84	<b>0.99</b>	0.95	<b>0.99</b>	0.54	<b>0.99</b>	<b>0.96</b>	<b>0.99</b>
purine	0.36	<b>0.99</b>	0.36	<b>0.99</b>	0.32	<b>0.99</b>	<b>0.37</b>	<b>0.99</b>
ylbH	0.01	0.5	0.01	<b>1.00</b>	<b>0.02</b>	<b>1.00</b>	0.01	0.33
cobalamin	None		0.84	0.82	0.72	<b>1.00</b>	<b>0.86</b>	0.99
lysine	0.79	<b>1.00</b>	<b>0.84</b>	0.82	0.72	<b>1.00</b>	0.74	<b>1.00</b>
SRP	0.1	0.99	0.1	<b>1.0</b>	<b>0.84</b>	0.98	0.77	0.98
RNaseP	<b>1.0</b>	<b>0.99</b>	<b>1.0</b>	<b>0.99</b>	<b>1.0</b>	<b>0.99</b>	<b>1.0</b>	<b>0.99</b>
FMN	<b>0.96</b>	<b>0.99</b>	<b>0.96</b>	<b>0.99</b>	<b>0.96</b>	<b>0.99</b>	<b>0.96</b>	0.98
glycine	None		0.08	0.98	None		<b>0.86</b>	<b>0.99</b>
preQ1	<b>0.01</b>	<b>0.04</b>	None		<b>0.01</b>	0.02	None	
ydaO	<b>0.97</b>	<b>1.00</b>	0.96	<b>1.00</b>	0.96	<b>1.00</b>	0.96	0.99
yybP	<b>0.26</b>	<b>1.00</b>	0.11	<b>1.00</b>	0.11	<b>1.00</b>	0.22	0.99
6S	0.09	<b>1.00</b>	<b>0.42</b>	0.92	0.09	<b>1.00</b>	0.29	<b>1.00</b>
ykoK	<b>0.96</b>	<b>1.00</b>	<b>0.96</b>	<b>1.00</b>	0.90	<b>1.00</b>	<b>0.96</b>	<b>1.00</b>
glmS	<b>0.95</b>	<b>1.00</b>	0.93	<b>1.00</b>	0.91	<b>1.00</b>	<b>0.95</b>	<b>1.00</b>
ykkC	None		None		<b>0.69</b>	<b>1.00</b>	<b>0.69</b>	<b>1.00</b>
moco	None		None		None		None	
SMK	<b>0.08</b>	<b>0.67</b>	None		None		None	
Median	0.31	<b>0.99</b>	0.56	<b>0.99</b>	0.61	<b>0.99</b>	<b>0.76</b>	<b>0.99</b>

# Structure-Structure alignment

- Compare the distance of two known RNA secondary structures
- It is more useful nowadays as technologies for experimentally probing RNA secondary structures become mature
- Can be used to define distance for RNA structure clustering and defining RNA families based on conserved RNA secondary structures

# Algorithm for structure-structure alignment

**Procedure** *AlignRNA*

**begin**

**for** intervals  $(i_1, j_1)$ ,  $1 \leq i_1 < j_1 \leq m$

    and  $(i_2, j_2)$ ,  $1 \leq i_2 < j_2 \leq n$

(\* Assume that the intervals are examined in lexicographically increasing order of widths \*)

$$Align[i_1, j_1, i_2, j_2] = \max \begin{cases} Align[i_1, j_1 - 1, i_2, j_2] + \gamma(s[j_1], '-') \\ Align[i_1, j_1, i_2, j_2 - 1] + \gamma('-', t[j_2]) \\ Align[i_1, j_1 - 1, i_2, j_2 - 1] + \gamma(s[j_1], t[j_2]) \end{cases}$$

**if** there exist  $i_1 \leq k_1 < j_1, i_2 \leq k_2 < j_2$

  s.t.  $(k_1, j_1) \in S_1, (k_2, j_2) \in S_2$

$$Align[i_1, j_1, i_2, j_2] = \max \begin{cases} Align[i_1, j_1, i_2, j_2], \\ Align[i_1, k_1 - 1, i_2, k_2 - 1] + \\ Align[k_1 + 1, j_1 - 1, k_2 + 1, j_2 - 1] \\ + \delta(k_1, j_1, k_2, j_2) + \gamma(s[k_1], t[k_2]) + \gamma(s[j_1], t[j_2]) \end{cases}$$

**end**



# Complexity of Structure-Structure alignment algorithm

- $O(n^4)$ , because  $k_1$  and  $k_2$  can only take on constant values
- Zhong and Zhang further reduced the complexity to  $O(n^3)$

**Table 1 Comparison on running time of ERA, LocARNA, and RNAforester**

RNA family	length (bp)	num. pairs	ERA (sec)	LocARNA (sec)	ERA vs. LocARNA (fold)	RNAforester (sec)	ERA vs. RNAforester (fold)
tRNA	78	21	0.017	0.100	5.882	0.047	2.765
Gly riboswitch	105	22	0.015	0.277	18.46	0.162	10.80
U12 spliceosome	160	42	0.035	0.311	8.886	0.657	18.77
Phage_pRNA	244	43	0.124	0.647	5.218	6.935	55.93
tmRNA	367	64	0.929	22.45	24.16	225.4	242.6
biocoid_3UTR	549	155	4.898	170.3	34.77	13.99	2.856
snR86	1004	333	53.15	4862	91.48	5.579	-9.527*
Sacc_telomerase	1162	181	23.93	522.3	21.82	3697	154.5

\*ERA is slower than RNAforester when aligning snR86 RNA structures.

# Sequence-structure alignment

- Compared with structure-structure alignment, we lose the information regarding the secondary structure of one of the sequences
- Naively, we need to try all possible, say  $k_2$ , for branching in the recursive function
- This naïve way gives rise to an  $O(n^5)$  algorithm

# Sequence-structure alignment

- However, we notice that we know the secondary structure of one RNA sequence.
- We also observe that the branching case is only used for considering multi-branch loop
- We only compute the expensive recursive function when we know that we are handling a multi-branch loop case (by looking at the know secondary structure)
- This will reduce the complexity to  $O(n^4)$
- Zhang et al. further reduced the complexity to  $O(n^3)$ , See Zhang et al. 2005 IEEE/ACM TCBB

# Sequence-structure alignment algorithm

**Procedure** *InferStructure*

**begin**

**for** intervals  $(i_1, j_1)$ ,  $1 \leq i_1 < j_1 \leq n$

    and intervals  $(i_2, j_2)$ ,  $1 \leq i_2 < j_2 \leq m$

(\* Assume that the intervals are examined in lexicographically increasing order of widths\*)

$$Align[i_1, j_1, i_2, j_2] = \max \begin{cases} Align[i_1 + 1, j_1, i_2, j_2] + \gamma(s[i_1], '-') \\ Align[i_1, j_1, i_2 + 1, j_2] + \gamma('-', t[i_2]) \\ Align[i_1 + 1, j_1, i_2 + 1, j_2] + \gamma(s[i_1], t[i_2]) \\ Align[i_1, j_1 - 1, i_2, j_2] + \gamma(s[j_1], '-') \\ Align[i_1, j_1, i_2, j_2 - 1] + \gamma('-', t[j_2]) \\ Align[i_1, j_1 - 1, i_2, j_2 - 1] + \gamma(s[j_1], t[j_2]) \end{cases}$$

**if**  $(i_1, j_1) \in S$  and

$t[i_2]$  and  $t[j_2]$  are complementary base-pairs

$$Align[i_1, j_1, i_2, j_2] = \max \begin{cases} Align[i_1, j_1, i_2, j_2], \\ \delta(i_1, j_1, i_2, j_2) + \gamma(s[i_1], t[i_2]) \\ \quad + \gamma(s[j_1], t[j_2]) + Align[i_1 + 1, j_1 - 1, i_2 + 1, j_2 - 1] \end{cases}$$

**else if**  $(i_1, j_1) \in S^c - S$  and

$(k, j_1) = \text{rightchild}(i_1, j_1)$

$$Align[i_1, j_1, i_2, j_2] = \max \begin{cases} Align[i_1, j_1, i_2, j_2], \\ \max_{i_2 < l < j_2} \{ Align[i_1, k - 1, i_2, l - 1] + Align[k, j_1, l, j_2] \} \end{cases}$$

**end**

# Simultaneous alignment and folding

- What if we know none of the two RNA sequences' secondary structures, and we want to compute the consensus structure that can be formed by them?
- Scoring: a combination of sequence & structural similarity and thermodynamic stability of the consensus structure
- Algorithm with only similarity: Simple! Just try all  $k_1$  and  $k_2$ !
- $O(n^6)$  time complexity; unfortunately, no easy way to reduce it

# Simultaneous alignment and folding

- Due to David Sankoff; also called Sankoff's algorithm
- Incorporating similarity and stability
- Under Zuker-Sankoff energy model

# Sankoff's algorithm

$$\begin{aligned}
 & C(i_1, j_1; i_2, j_2) \\
 = \min & \left\{ \begin{aligned}
 & e(s_1) + e(s_2) + D(i_1 + 1, j_1 - 1; i_2 + 1, j_2 - 1), \quad s_1, s_2 \text{ hairpins closed by} \\
 & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad (i_1, j_1), (i_2, j_2) \text{ respectively,} \\
 & \min \{ e(s_1) + e(s_2) + C(p_1, q_1; p_2, q_2) \\
 & \qquad \qquad \qquad + D(i_1 + 1, p_1; i_2 + 1, p_2) + D(q_1, j_1 - 1; q_2, j_2 - 1) \}, \\
 & \qquad \qquad \qquad s_1, s_2 \text{ are 2-loops closed by } (i_1, j_1), \\
 & \qquad \qquad \qquad (i_2, j_2) \text{ with } (p_1, q_1), (p_2, q_2) \text{ accessible,} \\
 & \qquad \qquad \qquad p_1 - i_1 + j_1 - q_1 - 2 \leq U, p_2 - i_2 + j_2 - q_2 - 2 \leq U, \\
 & \qquad \qquad \qquad \text{or one of } \begin{cases} s_1 = \phi & \text{and } (p_1, q_1) = (i_1, j_1) \\ s_2 = \phi & \text{and } (p_2, q_2) = (i_2, j_2), \end{cases} \\
 & \min_{\substack{i_1 < h_1 < j_1 - 1 \\ i_2 < h_2 < j_2 - 1}} \{ G(i_1 + 1, h_1; i_2 + 1, h_2) \\
 & \qquad \qquad \qquad + G(h_1 + 1, j_1 - 1; h_2 + 1, j_2 - 1) + 2A \},
 \end{aligned} \right.
 \end{aligned}$$

# Sankoff's algorithm cont.

$$G(i_1, j_1; i_2, j_2)$$

$$= \min \left\{ \begin{array}{l} C(i_1, j_1; i_2, j_2) + 2P + D(i_1, i_1; i_2, i_2) + D(j_1, j_1; j_2, j_2), \\ \min_{\substack{i_1 < h_1 < j_1 \\ i_2 < h_2 < j_2}} \min \left\{ \begin{array}{l} G(i_1, h_1; i_2, h_2) + (j_1 - h_1 + j_2 - h_2)Q \\ \quad + D(h_1 + 1, j_1; h_2 + 1, j_2), \\ G(i_1, h_1; i_2, h_2) + G(h_1 + 1, j_1; h_2 + 1, j_2), \\ (h_1 - i_1 + 1 + h_2 - i_2 + 1)Q \\ \quad + G(h_1 + 1, j_1; h_2 + 1, j_2) + D(i_1, h_1; i_2, h_2), \end{array} \right. \end{array} \right.$$

$$F(i_1, j_1; i_2, j_2) = \min \left\{ \begin{array}{l} C(i_1, j_1; i_2, j_2) + D(i_1, i_1; i_2, i_2) + D(j_1, j_1; j_2, j_2), \\ \min_{\substack{i_1 \leq h_1 < j_1 \\ i_2 \leq h_2 < j_2}} \{F(i_1, h_1; i_2, h_2) + F(h_1 + 1, j_1; h_2 + 1, j_2)\}, \\ D(i_1, j_1; i_2, j_2), \end{array} \right.$$