

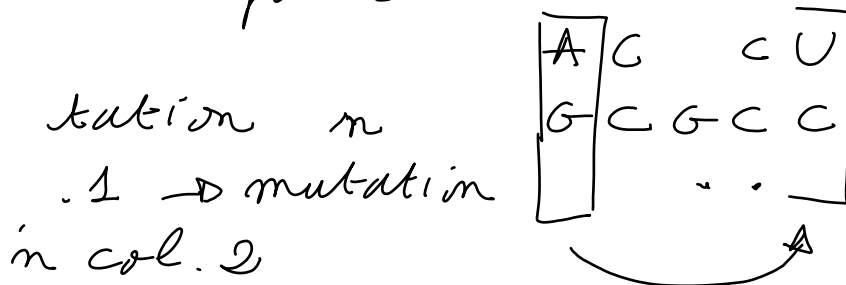
# NA sequence structure alignment

## ⊖ Motivation

- Single structure prediction is limited (~70% base pairs correctly predicted)

- structures predictions are obtained through comparative modeling.

1. obtain a sequence alignment of homologous NA  
bind minimal set of overlapping pairs

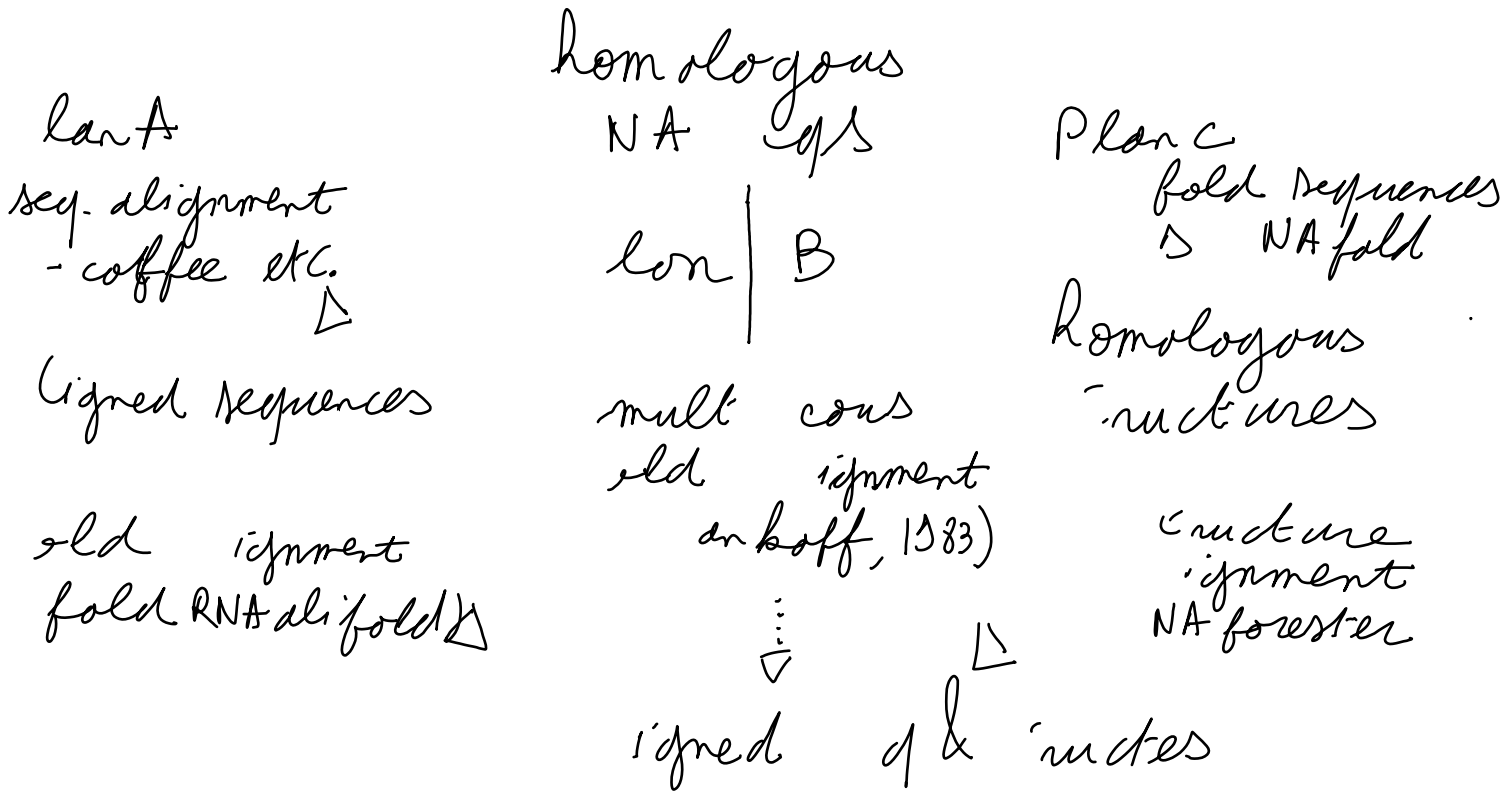


Ex: 6S RNA at base from Gutell

Very accurate → 9% accurate

conclusion: could be used in next in. how? (information content have an alignment)

# Strategies



Problems:

Seq A may fail if sequences are almost identical (0%)

Seq C structures may not be available  
alignments are not accurate

Solution: Don't do things sequentially  
simultaneous - folding & alignment  
Sankoff (1983)

question, how to do SFA?

○ sequence alignment

dynamic time & space with edit distance

$$D(i, j) = \max \begin{cases} D(i, j-1) + \delta(-, w'_j) \\ D(i-1, j) + \delta(w_i, -) \\ D(i-1, j-1) + \delta(w_i, w'_j) \end{cases}$$

$$D(i, 0) = D(0, i) \geq \delta(w_i, -) \quad \sum_i \delta(-, w_i)$$

$\delta(-, x) = -\text{gap penalty}$

$$\delta(x, y) = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{otherwise} \end{cases}$$

sequences  $O(n^2)$  time (full  $\Delta$  can be optimized)

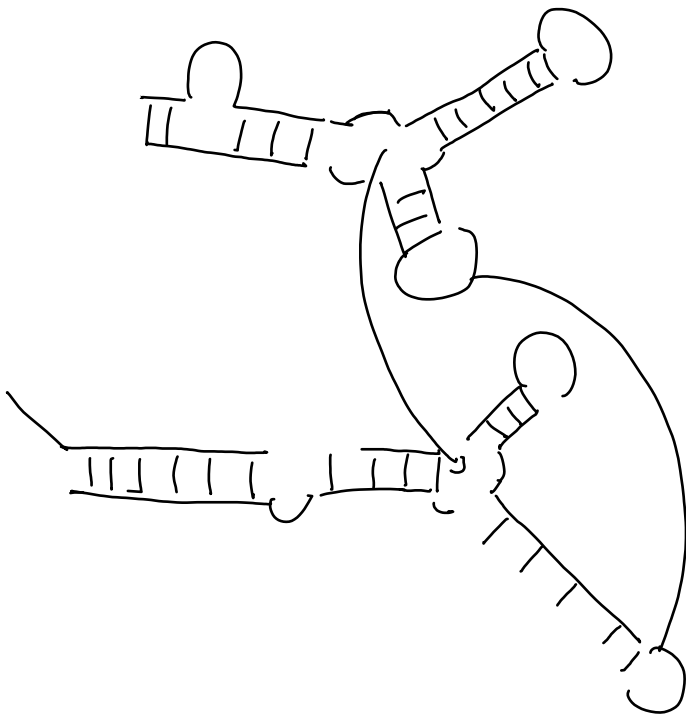
○ folding

$(n^3)$  time  $O(n^2)$  space (NJ, Zuker, etc.)

▷ we should be to do both partition in dynamic time & space.

# ○ Sankoff - algorithm

B.1 concept: compatible structures



similar branching structures.

related to constraint assignments

1.  $\exists (i, j) \in S$  s.t.  $i < \pi < j$   
 $\nexists (i', j')$  s.t.  $i < i' < \pi < j' < j$   
 $\rightarrow$  is accessible from  $(i, j)$

an structure alignment is .t.

- accessible union of external & multi loop are ignored
- incoming pairs of regions are aligned
- ops index & aligned, inserted, deleted

①  $i_1 < \dots < i_k$  positions of an external ~~loop~~ <sup>pair</sup>  
 . accessible <sup>pair</sup> positions of a  
 multi-loop  $n$   $w$   
 $i'_1 < \dots < i'_l$  on  $w'$

→ the algo produces an alignment  
 such that  
 $k = l$   
 •  $(i_f, i'_g) \in -$  iff  $(i_f, i'_g) \in B$

same branching figuration.

recursive equations

$$D(i, j; i', j') = \max \left\{ \begin{array}{l} D(i, j; i', j'-1) + \delta(w_j, w'_{j-1}) \\ D(i, j-1; i', j') + \delta(w_{j-1}, w'_j) \\ D(i, j; i', j') + \delta(w_j, w'_j) \end{array} \right.$$

$$\delta(x, y) = \begin{cases} (w_j, w'_j) & \text{if } x == y \\ x & \text{otherwise} \end{cases}$$

→ sequence alignment of loop regions.

## 2. anhoff Algorithm

$D[(i_1, i_2), (j_1, j_2)] \leftarrow$  best alignment score (sequence only)  
 $\tau - i_1 j_1$  &  $B_{i_2 j_2}$

$V[(i_1, i_2), (j_1, j_2)] \leftarrow$  min - alignment (seq + struct.)  
 $- - i_1 j_1$  &  $B_{i_2 j_2}$

$\downarrow \leftarrow$  same with multi-loop  
 nality for internal bases

$[ (i_1, i_2), (j_1, j_2) ]$  — min cost alignment (seq & struct.)  
 of  $A_{i_1 j_1}$  &  $B_{i_2 j_2}$  st.  $(i_1, j_1)$  &  $(i_2, j_2)$   
 e base pairs

equations:

$D$  can be filled using eddleman-Wunsch

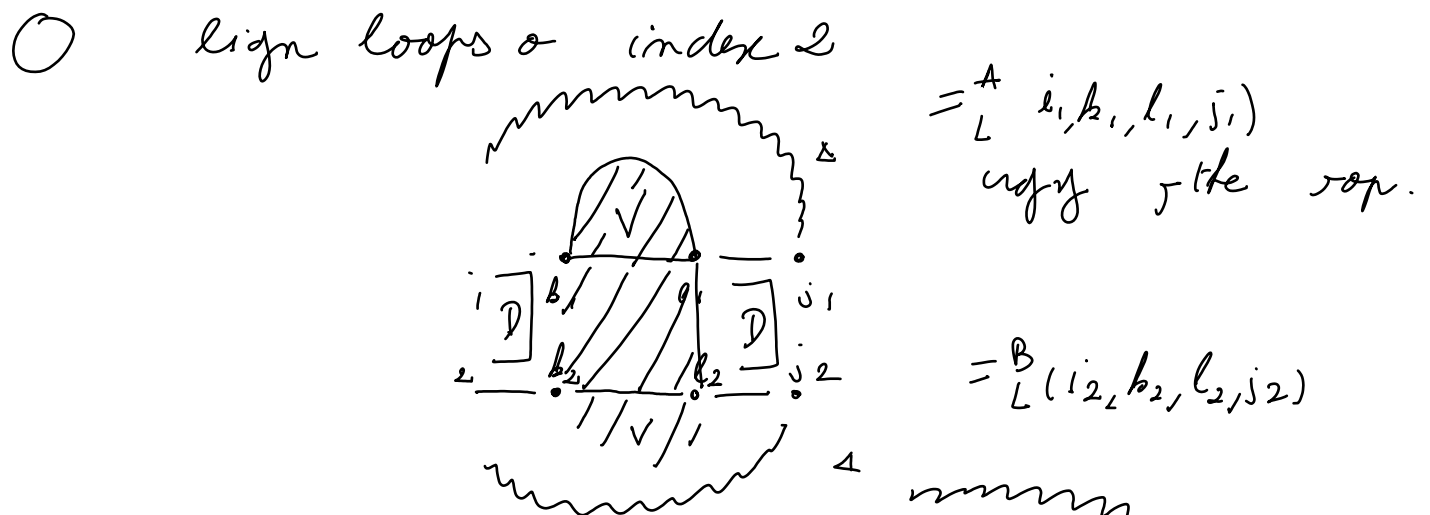
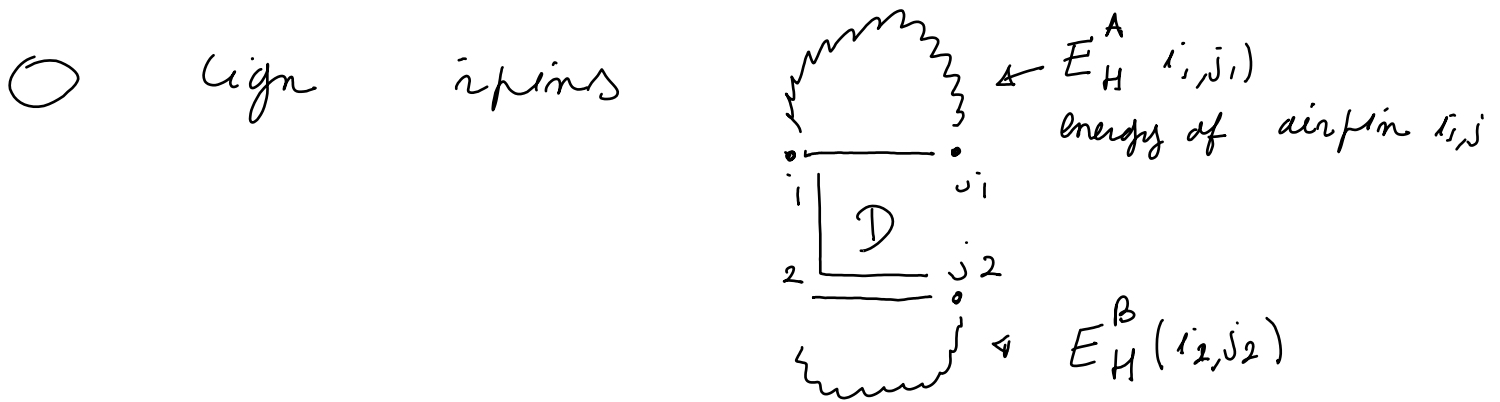
$V[(i_1, i_2), (j_1, j_2)]$  / base pair at  
 extremities

$$\bullet W[(i_1, i_2), (j_1, j_2)] = \min_{k_1, k_2} V[(i_1, i_2), (k_1, k_2)] + W[(k_1+1, k_2+1), (j_1, j_2)]$$

/ all decomp.

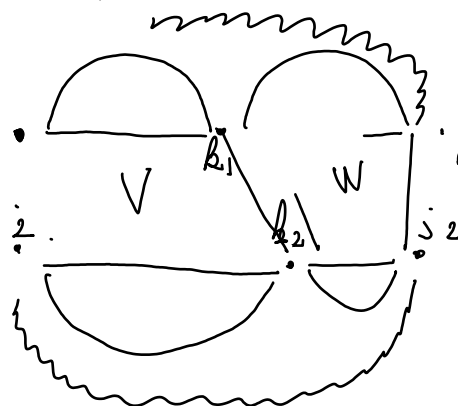
$D[(i_1, i_2), (j_1, j_2)]$  // no structure

$$\begin{aligned}
 & D[(i_1, j_1), (i_2, j_2)], E_H^A(i_1, j_1) + E_H^B(i_2, j_2) \\
 & \min_{l_1, l_2} = E_L^A(i_1, b_1, l_1, j_1) + E_L^B(i_2, b_2, l_2, j_2) \\
 & + D[(i_1, j_1), (b_1, l_1)] + D[(b_1, l_1), (b_2, l_2)] \\
 & + V[(b_1, l_1), (b_2, l_2)] \\
 & =_{\text{close}} E_H^A(i_1, j_1) + E_H^B(i_2, j_2) \\
 & \min_{l_1, l_2} V[(i_1, j_1), (b_1, l_1)] + W[(b_1, l_1), (b_2, l_2)] + W[(b_2, l_2), (i_2, j_2)]
 \end{aligned}$$



○ sign branching

Rd: we need to concatenate  
 2x W or ensure  
 3 rem



• solution of the problem is found in  
 $N \left[ \begin{pmatrix} 1 \\ 1 \end{pmatrix} \begin{pmatrix} |A| \\ |B| \end{pmatrix} \right]$  & backtracks

complexity  $O(|A|^3 |B|^3)$  - time  
 $(|A| |B|^2)$  - space

▷ too high for immediate application  
especially in 1985!

### 3.4. Discussion

Why does it work? FE of single seq is not accurate, but if a structure has a low energy & homologous sequences simultaneously then this is m - likely a better prediction

○ anhoff-style algo in practice

first practical implementations appeared in the 2000's

Dynalign in RNA structure) Mathews Turner 2002)

time complexity but efficient implementation if restricted. (See next page)



Dynalign complexity  $O((\min(|A|, |B|))^3 K^3)$

Here  $K$  is the maximum distance between aligned nucleotides.

⇒ allows to restrict search depth

Application: detection of ncRNA  
(Uzilov - al., 200)

not a cal:

- u-c-d-base of ncRNA
- align query RNA with seq in DB
- generate model of for this var of RNA & calculate Z-score
- accept as cRNA if Z-score is significant

M np Hofacker et al., 2004)

Dynalign or another implementation  $\left\{ \begin{array}{l} \text{Sankoff} \\ \text{Lyo (foldalign, Orokin et al. 1997)} \end{array} \right.$   
or not implement the full model

idea: se base pair proba calculated by  
Cashill, 1990) or sped up Sankoff.

you?

• in the algo only  
base pair with proba higher than  
threshold.

benefits:

- deals with a number of potential base pair lined with the length of RNAs deals with base pairs (i.e. not stacks) → simpler

methodology can be used to align stochastic contact matrices. (re informative)  
(can be tuned with convex opt.)

Application: LocaRNA (Will et al. 200) (re implementation)

- clustering - ncRNA.
  - 1 take a set of ncRNA
  - 2 align all sets of ncRNA and compute a score (distance)
  - 3 cluster RNA used on this distance matrix (here weighted pair method algo)  
instances  $d(i,j) = \begin{cases} 0, & \text{score} \end{cases}$   
here  $q$  is the 9%-quantile of all pairwise scores.

is method used on Rfam (ds)

- discover some classes of ncRNA

- found classes of ncRNA in *iona intestinalis* genome