Foldalign 2.2: multithreaded structural RNA alignment

Daniel Sundfeld^{1,3} Jakob Havgaard^{1,2} Alba C. M. A. de Melo³ Jan Gorodkin^{1,2}

¹Center for non-coding RNA in Technology and Health

²IKVH University of Copenhagen

³Department of Computer Science University of Brasilia

19/02/2015 - Winterseminar

Foldalign





Algorithms for simultaneously aligning and folding:

- the probabilistic methods: usually based on Stochastic Context Free Grammars, that use statistical learning to creates a solution based on previously known answers;
- **Sankoff-like methods**: use some method of minimization of the free energy to calculates how elements of an RNA structure contribute to free energy.

*M. Ziv-Ukelson, I. Gat-Viks, Y. Wexler, and R. Shamir. A faster algorithm for rna co-folding. In Proceedings of the 8th International Workshop on Algorithms in Bioinformatics, pages 174–185, 2008.

The time complexity $O(L^6)$, where L is the sequence length of the sequences, makes the Sankoff algorithm prohibited for long sequences. Foldalign use several constraints to reduce the time requirement:

• The final alignment is limited to λ nucleotides, and the maximum length difference between any two subsequences being aligned is limited to δ nucleotides;

• Subalignments with score bellow to a threshold are eliminated. With this constraints, the time complexity is reduced to $O(L^2\lambda^2\delta^2)$ The simplified Foldalign recursion:

$$\begin{pmatrix} D_{i+1,j-1,k+1,l-1} + S_{bp}(n_i, n_j, n_k, n_l, \sigma_{bp}) & (a) \\ D_{i+1,j-1,k,l} + S_{bpil}(n_i, n_j, -, -, \sigma_{bpil}) & (b) \\ D_{i,j,k+1,l-1} + S_{bpik}(-, -, n_k, n_l, \sigma_{bpik}) & (c) \\ \end{pmatrix}$$

$$D_{i+1,j-1,k,l} + S_{bpil}(n_i, n_j, -, -, \sigma_{bpil})$$
 (b)

$$\mathcal{D}_{i,j,k+1,l-1} + \mathcal{S}_{bpiK}(-,-,n_k,n_l,\sigma_{bpiK})$$
 (c)

$$D_{i+1,j,k+1,l} + S_{al}(n_i, n_k, \sigma_{al})$$
(d)

$$D_{i,j-1,k,l-1} + S_{ar}(n_j, n_l, \sigma_{ar})$$
(e)

$$D_{i,j,k,l} = \max \begin{cases} D_{i,j-1,k,l-1} + S_{gll}(n_i, -, \sigma_{gll}) & (f) \\ D_{i+1,j,k,l} + S_{gll}(n_j, -, \sigma_{glrl}) & (g) \\ D_{i,j-1,k,l} + S_{glk}(n_j, -, \sigma_{glrl}) & (g) \\ D_{i,j,k+1,l} + S_{glK}(-, n_k, \sigma_{glK}) & (h) \\ D_{i,j,k,l-1} + S_{grK}(-, n_l, \sigma_{grK}) & (i) \end{cases}$$

$$O_{i,j-1,k,l} + S_{glrl}(n_j, -, \sigma_{glrl})$$
 (g)

$$\mathcal{D}_{i,j,k+1,l} + \mathcal{S}_{glK}(-, n_k, \sigma_{glK})$$
 (h)

$$D_{i,j,k,l-1} + S_{grK}(-, n_l, \sigma_{grK})$$
 (i)

Image: Image:

$$\sum_{\substack{i < m < j \\ k < n < l}}^{i < m < j} \{ D_{i,m,k,n} + D_{m+1,j,n+1,l} + \sigma_{mbl} \}$$
(j)

(1)

Foldalign Graphical Representantion



J. Gorodkin and W. L. Ruzzo. RNA Sequence, Structure and Function: Computational and Bioinformatic Methods. Humana Press, 2014.

3

• • • • • • • • • • • •

1:	for $i = Length_I ightarrow 1$ do
2:	for $k = Length_{\mathcal{K}} ightarrow 1$ do
3:	for $j=1 ightarrow\lambda$ do
4:	for $l=(j-\delta) ightarrow (j+\delta)$ do
5:	if $Ms_{i,j,k,l}$ can be the right part of bifurcation then
6:	$MI_{i,j,k,l} = Ms_{i,j,k,l}$
7:	end if
8:	$E_{xpand_Position(i, j, k, l)}$
9:	for $W_m = 1 ightarrow (\lambda - j)$ do
10:	for $W_n = (W_m - \delta) ightarrow (W_m + \delta)$ do
11:	$E_{xpand_MB(i, j, k, l, W_m, W_n)}$
12:	end for
13:	end for
14:	end for
15:	end for
16:	end for
17:	end for

・ロト ・聞ト ・ヨト ・ヨト

The multithreading algorithm is based on the fact that cells in the i loop might be calculated in parallel. Some data dependency exist inside the loop, conditions are added to avoid race conditions.

 $N_{threads}$ are created, each one have different *i* values:

Thread 1: Thread 2: (...) 1: for $i = L_I \rightarrow 1$ step $N_{threads}$ do 1: for $i = (L_I - 1) \rightarrow 1$ step $N_{threads}$ do 2: for $k = L_K \rightarrow 1$ do 2: for $k = L_K \rightarrow 1$ do 3: 3: (...) (...) 16: end for 16: end for 17: end for 17: end for

- E > - E >

- 1: for i = 20 do 2: for k = 20 do 3: (...)
- 16: end for 17: end for

- 1: for i = 20 do 2: for k = 19 do 3: (...)
- 16: end for 17: end for

1: for i = 20 do 2: for k = ... do 3: (...)

16: end for 17: end for

- 1: for i = 20 do 2: for k = 1 do 3:
- (...)

16: end for 17: end for

- 1: for i = 19 do 2: for k = 20 do 3: (...)
- 16: end for 17: end for

- 1: for i = 19 do 2: for k = 19 do 3: (...)
- 5. (...)
- 16: end for 17: end for

1: for i = 20 step 2 do 2: for k = 20 do 3: (...)

16: end for 17: end for Thread 2:

1: for *i* = 19 step 2 do 2: for *k* = waiting do 3: (...)

16: end for 17: end for

1: for i = 20 step 2 do 2: for k = 19 do 3: (...)

16: end for 17: end for Thread 2:

1: for *i* = 19 step 2 do 2: for *k* = waiting do 3: (...)

16: end for 17: end for

1: for i = 20 step 2 do 2: for k = 18 do 3: (...)

16: end for 17: end for Thread 2:

1: for i = 19 step 2 do 2: for k = 20 do 3: (...)

16: end for 17: end for

1: for i = 20 step 2 do 2: for k = (...) do 3: (...)

16: end for 17: end for Thread 2:

1: for i = 19 step 2 do 2: for k = (...) do 3: (...)

16: end for 17: end for

1: for i = 18 step 2 do 2: for k = 20 do 3: (...)

16: end for 17: end for Thread 2:

1: for *i* = 17 step 2 do 2: for *k* = waiting do 3: (...)

16: end for 17: end for

1: for i = 18 step 2 do 2: for k = 19 do 3: (...)

16: end for 17: end for Thread 2:

1: for *i* = 17 step 2 do 2: for *k* = waiting do 3: (...)

16: end for 17: end for

1: for i = 18 step 2 do 2: for k = 18 do 3: (...)

16: end for 17: end for Thread 2:

1: for i = 17 step 2 do 2: for k = 20 do 3: (...)

16: end for 17: end for

Other mechanisms are necessary to the final implementation:

- **Threadpoll:** To avoid thread creation overhead, a fixed number of threads are created and each thread wait for a *i* value.
- **Protection:** The global memory need protection mechanisms to avoid simultaneous writes.
- and other implementation details...

Results Multithreaded Programming: Expectation vs. Reality



Results

Random Sequences Performance



2000 nucleotides







Sundfeld et al

Foldalign 2.2

э 19/02/2015 12 / 16

Results

Real Sequences Performance

16S Ribosomal RNA from the Gutell website http://www.rna.icmb.utexas.edu/ $(|S_1| = 2315, |S_2| = 2019, \lambda = 2019, \delta = 100).$



< A

3 🕨 🖌 3

A multithreaded version of Foldalign is now available.

The implementation were carefully designed to keep all previous functions, and it open new possibilities to search for structured RNAs in much longer sequences in reasonable time, and in our experiments, it was possible to search 2.5-2.7 times faster with random sequences, and 4.18 times faster with 8 cores, 5.03 times faster with 14 cores.

Acknowledgements









イロト イ団ト イヨト イヨト

 →
 ≡
 つ<</th>

 19/02/2015
 15 / 16

Thank you! Questions?

-

イロン イヨン イヨン イ