# Lifting the Prediction of RNA Pseudoknots to their Alignment
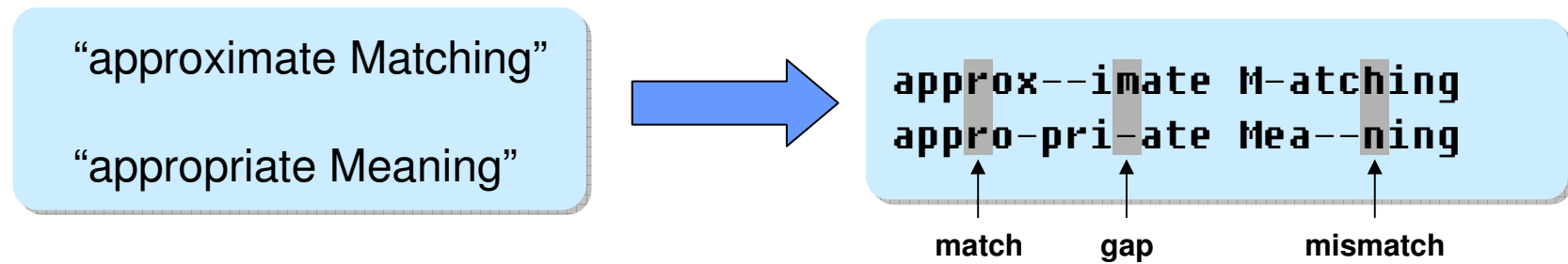
by Mathias Möhl
Programming Systems Lab., Saarland University, Prof. Smolka

joint work with Rolf Backofen and Sebastian Will
Chair for Bioinformatics, Albert-Ludwigs-Universität Freiburg

# Alignment

**sequence alignment:**   comparison of sequences (RNAs, DNAs, Proteins)

> "approximate Matching"
>
> "appropriate Meaning"

```
approx--imate M-atching
appro-pri-ate Mea--ning
```

match        gap        mismatch

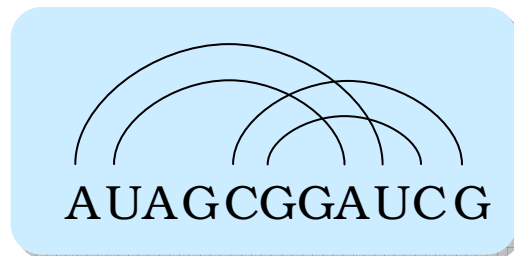**structure alignment:** also consider secondary structure in comparison

```
..((( (((((((... ..........))))))) ))).... 
ACGUG-UGAUGGGAGG-UACAAGCAACCCCAUUA-CAUAUUA
UGCUGCU-AUAAUA-GAUAGA-GA-AGGUUAU-AGCAGACUA
..((((( (((((. ...... .. ..))))) )))))....
```
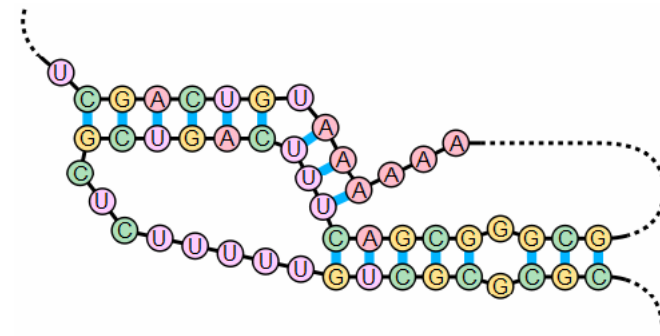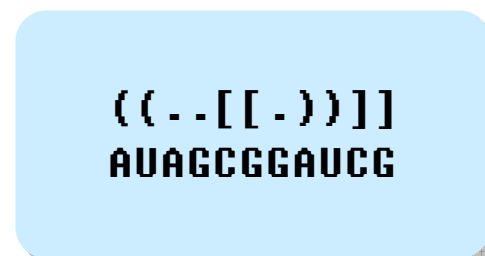
# Pseudoknots

an RNA structure contains pseudoknots

...if it contains any crossing arcs



AUAGCGGAUCG
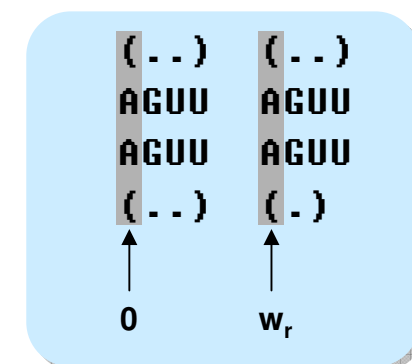
structures *without* pseudoknots can easily be decomposed

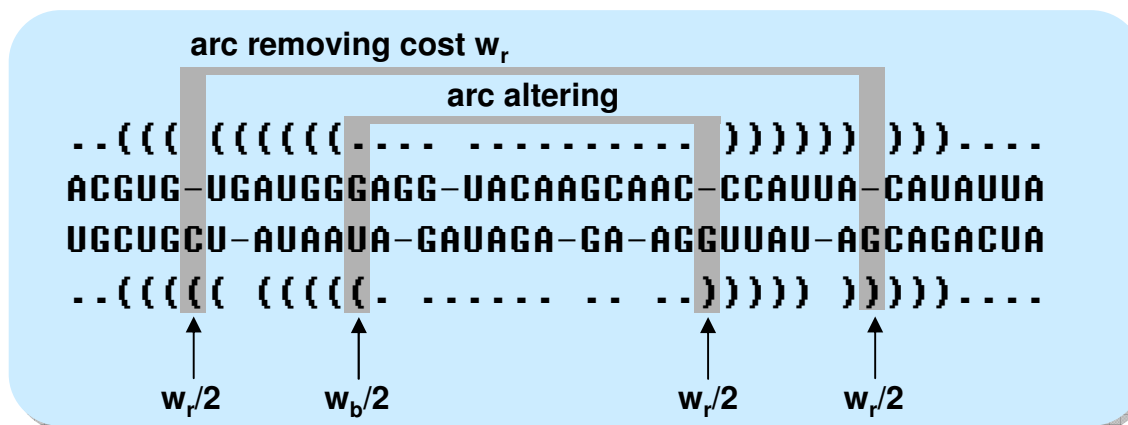→efficient dynamic programming algorithms

more difficult for pseudoknots

...if it is not well-bracketed

```
((..[[.))]]
AUAGCGGAUCG
```

# Optimal Alignment

optimal alignment ≈ minimum edit distance

| edit operation | cost | example |
|---|---|---|
| base deletion | $w_d$ | $\overset{.}{A} \rightarrow -$ |
| base substitution | $w_m$ | $\overset{.}{A} \rightarrow \overset{.}{C}$ |
| arc mismatch | $w_{am}/2$ per mismatched base | $\overset{(}{A} \quad \overset{)}{U} \rightarrow \overset{(}{C} \quad \overset{)}{G}$ |
| arc breaking | $w_b$ | $\overset{(}{A} \quad \overset{)}{U} \rightarrow \overset{.}{A} \quad \overset{.}{U}$ |
| arc removing | $w_r$ | $\overset{(}{A} \quad \overset{)}{U} \rightarrow - \quad -$ |
| arc altering | $w_a = (w_b + w_r)/2$ | $\overset{(}{A} \quad \overset{)}{U} \rightarrow \overset{.}{A} \quad -$ |

arc removing cost $w_r$

arc altering

```
..((((  ((((((.---  ...........  ))))))  )))....
ACGUG-UGAUGGGAGG-UACAAGCAAC-CCAUUA-CAUAUUA
UGCUGCU-AUAAUA-GAUAGA-GA-AGGUUAU-AGCAGACUA
..(((((( (((((.  .......  ..  ..)))))  ))))))....
```

$w_r/2 \qquad w_b/2 \qquad\qquad w_r/2 \qquad w_r/2$

```
(..)    (..)
AGUU    AGUU
AGUU    AGUU
(..)    (.)
```

$0 \qquad w_r$

# Optimal Alignment

optimal alignment ≈ minimum edit distance

| edit operation | cost | example |
|---|---|---|
| base deletion | w | |
| base subs | | |
| arc misma | | |
| arc breaki | | |
| arc remov | | |
| arc alterin | | |

cost of preserving one end of arc depends on other end

all other costs are independent

arc removing cost $w_r$

arc altering

```
..(((  ((((((.---  ..........  ))))))  )))....
ACGUG-UGAUGGGAGG-UACAAGCAAC-CCAUUA-CAUAUUA
UGCUGCU-AUAAUA-GAUAGA-GA-AGGUUAU-AGCAGACUA
..(((((  (((((.  .......  ..  ..)))))  )))))....
```

$w_r/2$     $w_b/2$          $w_r/2$     $w_r/2$

```
( .. )    ( .. )
AGUU      AGUU
AGUU      AGUU
( .. )    ( . )
```

0         $w_r$

# Structure Prediction vs. Alignment

**structure prediction**

**sequence structure alignment**

auagcagcuc

a u a g g g c u c

g c g c a g c u c

↓

↓

auagcagcuc

a u a g   g g c u c

g c   g c a g c u c

optimization problem:
minimum free energy structure

optimization problem:
minimum edit distance

both solvable with dynamic programming (DP)
where recursive decomposition is based on secondary structure
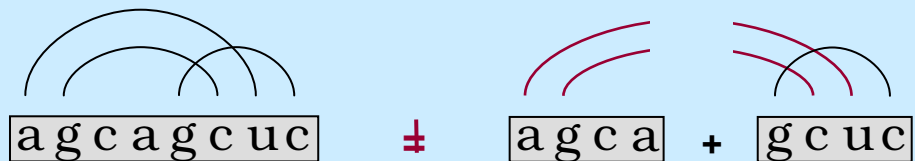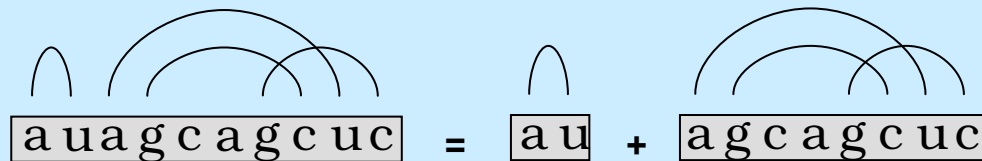
# What has been done before

| class of structures | pseudoknot free | R&G | A&U | L&P | D&P | R&E | arbitrary pseudoknots |
|---|---|---|---|---|---|---|---|
| **prediction** | [Zucker1981] [Sankoff1985] | [Reeder2004] | [Uemura1999] [Akutsu2000] | [Lyngso2000] | [Dirks2003] | [Rivas1999] | [Lyngso2000] |
| time | $O(n^3)$ | $O(n^4)$ | $O(n^4)$ | $O(n^5)$ | $O(n^5)$ | $O(n^6)$ | NP-complete |
| space | $O(n^2)$ | $O(n^2)$ | $O(n^3)$ | $O(n^3)$ | $O(n^4)$ | $O(n^4)$ | |
| **alignment** time | $O(n^2m^2)$ | | | | | $O(n^5m^5)$ | NP-complete |
| space | $O(nm)$ | | | | | $O(n^4m^4)$ | |
| | [Jiang2002] | | | | | [Evans2006] | [Evans1999] |
| | $O(mn^3)$ $O(mn^2)$ | $O(mn^4)$ $O(mn^2)$ | $O(mn^4)$ $O(mn^3)$ | $O(mn^5)$ $O(mn^3)$ | $O(mn^5)$ $O(mn^4)$ | $O(mn^6)$ $O(mn^4)$ | |
| | **new contribution** | | | | | | |

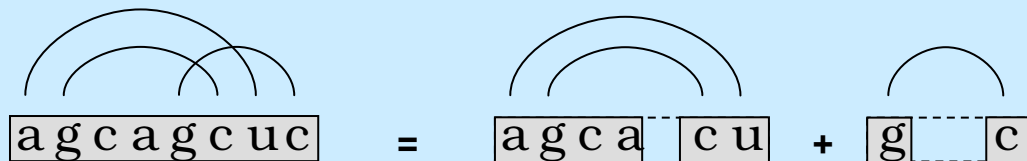**further alignment approaches:**

- fixed parameter tractable DP algorithms for arbitrary pseudoknots      [Evans1999],[Möhl2008]
- approaches based on integer linear programming (ILP)      [Lenhof1998],[Bauer2007],...

# RNA Structure Prediction



**General Framework** compose optimal structure recursively from *arc complete* fragments

a u a g c a g c u c = a u + a g c a g c u c

a g c a g c u c ⧧ a g c a + g c u c **NOT arc complete**

a g c a g c u c = a g c a c u + g c **arc complete**

- common base of all DP-based RNA alignment and structure prediction algorithms algorithms
- algorithms vary in
  - kinds of considered fragments
  - ways how to combine fragments

trade-off

complexity ↔ class of structures

# RNA structure prediction

**formal notion for that**

algorithms vary in

- kinds of considered fragments
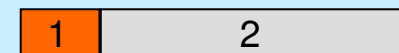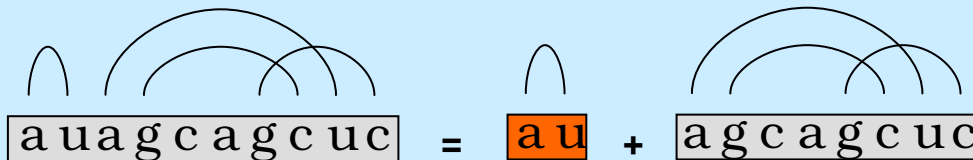- ways how to combine fragments

*number of intervals* of a fragment



1 interval

2 intervals
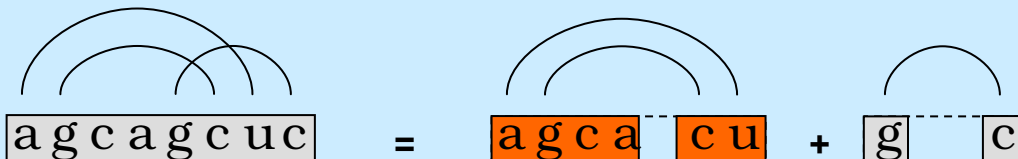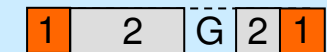
3 intervals

interval    fragment
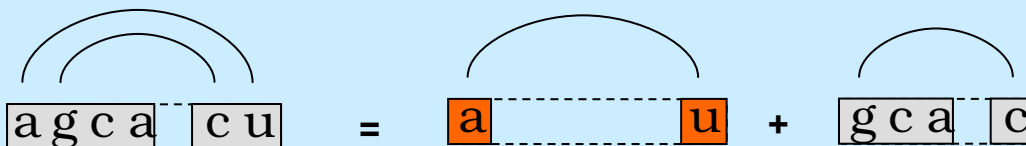
*split type*

$a\ u\ a\ g\ c\ a\ g\ c\ u\ c$ = $a\ u$ + $a\ g\ c\ a\ g\ c\ u\ c$

split type $T_1$ = 12

$a\ g\ c\ a\ g\ c\ u\ c$ = $a\ g\ c\ a\quad c\ u$ + $g\quad c$

split type $T_2$ = 1212

$a\ g\ c\ a\quad c\ u$ = $a\quad u$ + $g\ c\ a\quad c$

split type $T_3$ = 12G21

# A Structure Prediction Algorithm Scheme

optimal structure for fragments with 1 interval

= min

O(n²) space

all instances of split type $T_1$ — O(n³) time

all instances of split type $T_2$ — O(n⁴) time

⋮

all instances of split type $T_k$ — O(n⁵) time

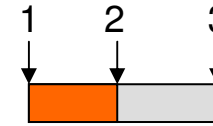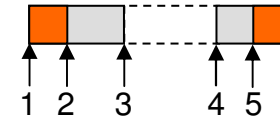1    2    3

optimal structure for fragment with 2 intervals

= min

O(n⁴) space

all instances of split type $T'_1$ — O(n⁴) time

⋮

all instances of split type $T'_{k'}$ — O(n⁶) time

1  2   3      4  5  6

**correctness**

correct for the class of structures that can be composed that way
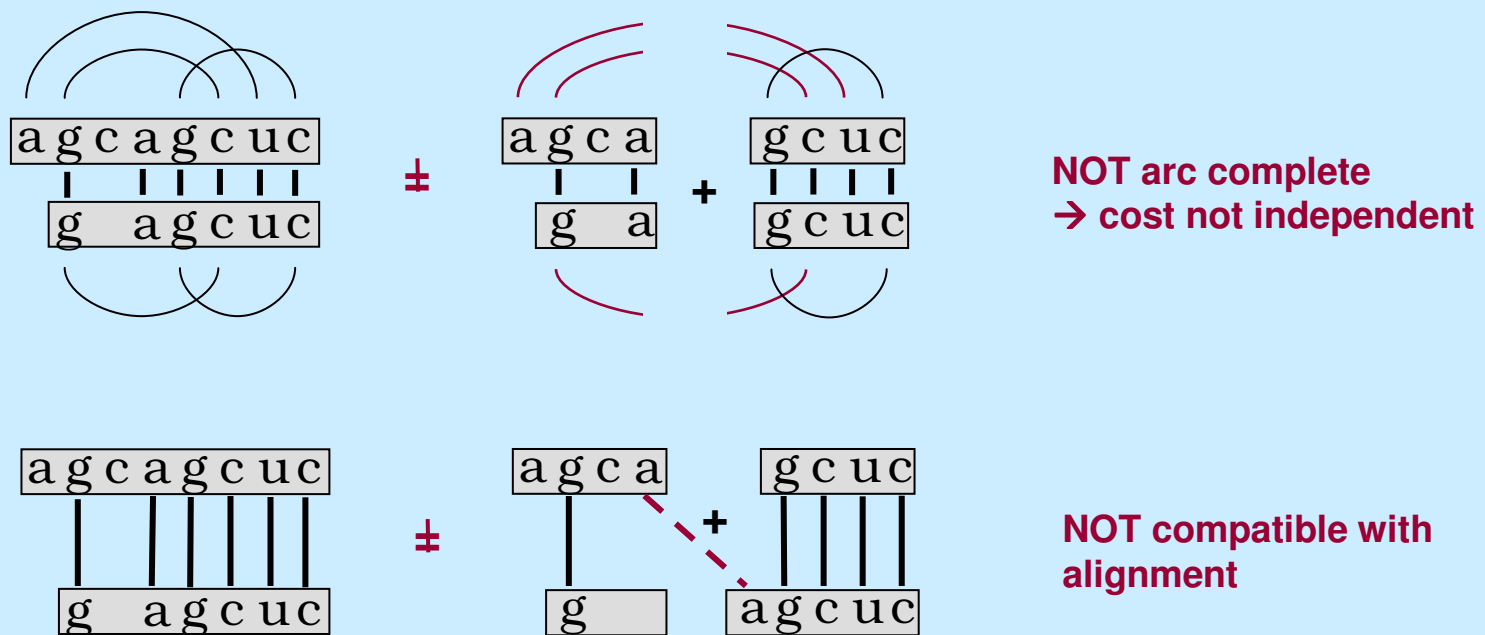(→ for which the case distinction is exhaustive)

**complexity**

time complexity:  depends on number of instances of considered split types
space complexity:  depends on degree of considered fragments

# Alignment

recursive split of alignments: **recursive split of 2 sequences simultaneously**

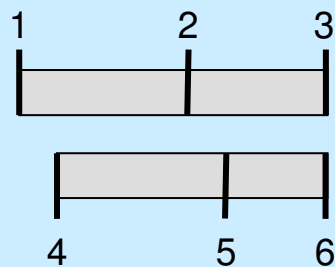split must be arc complete and compatible with alignment



**NOT arc complete**
**→ cost not independent**

**NOT compatible with alignment**

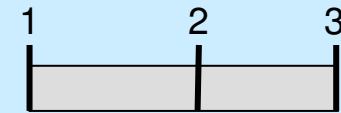**→ Idea:** search over all possible splits

# Alignment vs. Prediction

**Idea:** search over all possible splits

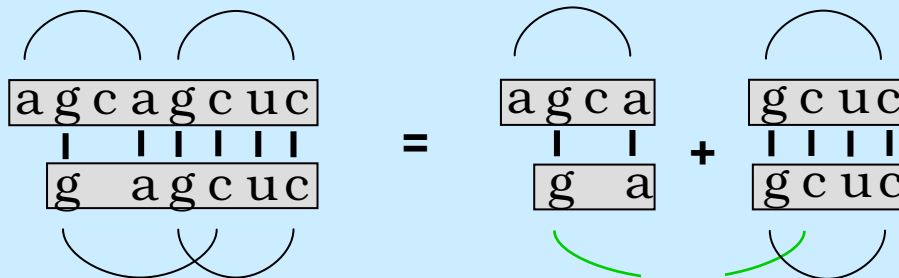**bad news:** there exist a lot more possibilities to split an alignment
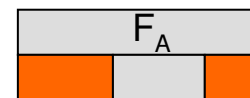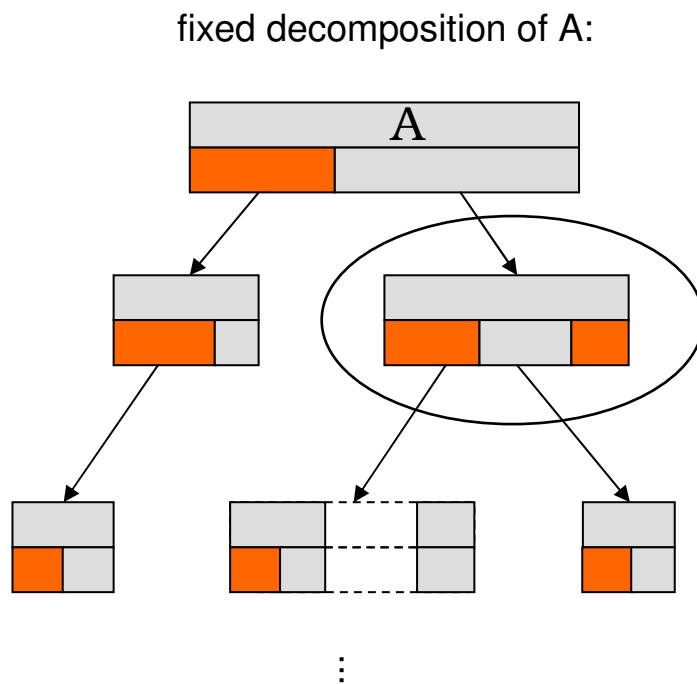


$O(n^{6)}$ instances

$O(n^{3)}$ instances

**good news:** breaking arcs in only one sequence is ok

# Alignment Algorithm

**idea** align RNAs A and B as follows:

consider a fixed decomposition of A , try to align it to all possible decompositions of B



fixed decomposition of A:

split type T

for all $F_B$ (i.e. for all i,j):

optimal alignment of

$$\frac{F_A}{F_B} = \min_{k,l}$$

fix

# Alignment Algorithm

for each decomposition of the first sequence there exists an appropriate decomposition of the second sequence



fixed decomposition of A:



split type T

for all $F_B$ (i.e. for all i,j):

optimal alignment of

$$\frac{F_A}{F_B} = \min_{k,l}$$

fix

i    j      i   k   l   j

# Alignment Algorithm

more general:

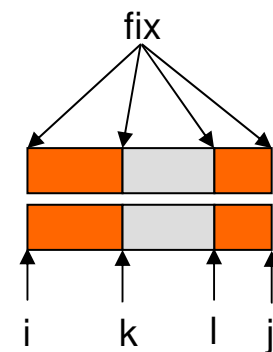for any fragments $F_A$ and $F_B$, and any arc complete split $(F_A^1, F_A^2)$ of $F_A$ with some type T

$$C(F_A, F_B) = \min_{T-\text{split } (F_B^1, F_B^2) \text{ of } F_B} \left\{ C(F_A^1, F_B^1) + C(F_A^2, F_B^2) \right\}$$

for all $F_B$ (i.e. for all i,j):

optimal alignment of

$F_A$

$F_B$ $= \min_{k,l}$

fix

i        j          i    k    l    j

⋮

# Alignment Algorithm Complexity

If RNA structure A has been predicted by some structure prediction algorithm X,
there exists a decomposition of A that uses only the split types used in algorithm X

# Extensions of the Scheme

| class of structures | | pseudoknot free | R&G | A&U | L&P | D&P | R&E | arbitrary pseudoknots |
|---|---|---|---|---|---|---|---|---|
| prediction | time | $O(n^3)$ | $O(n^4)$ | $O(n^4)$ | $O(n^5)$ | $O(n^5)$ | $O(n^6)$ | NP-complete |
| | space | $O(n^2)$ | $O(n^2)$ | $O(n^3)$ | $O(n^3)$ | $O(n^4)$ | $O(n^4)$ | |
| alignment | time | $O(n^2m^2)$ | | | | | $O(n^5m^5)$ | NP-complete |
| | space | $O(nm)$ | | | | | $O(n^4m^4)$ | |
| | | $O(mn^3)$ | $O(mn^4)$ | $O(mn^4)$ | $O(mn^5)$ | $O(mn^5)$ | $O(mn^6)$ | |
| | | $O(mn^2)$ | $O(mn^2)$ | $O(mn^3)$ | $O(mn^3)$ | $O(mn^4)$ | $O(mn^4)$ | |

**need extensions of the scheme**

more complex pseudoknots

- **extended split types**
  - need to capture exactly the splits considered by the algorithms
  - modify alignment computation to maintain correctness for these types
- **optimize space consumption**

→ **all tricks of the prediction algorithms can be transferred to alignment**

# Practical Evaluation: PKalign

COMPALIGN SCORE



- aligned pseudoknot structures of Rfam Database  (contains hand cured reference alignments)
- sequence length up to 125

# Conclusions

**RNA structure alignment and prediction**

- complexity depends on the structures

| class of structures | | pseudoknot free | R&G | A&U | L&P | D&P | R&E | arbitrary pseudoknots |
|---|---|---|---|---|---|---|---|---|
| **prediction** | time | $O(n^3)$ | $O(n^4)$ | $O(n^4)$ | $O(n^5)$ | $O(n^5)$ | $O(n^6)$ | NP-complete |
| | space | $O(n^2)$ | $O(n^2)$ | $O(n^3)$ | $O(n^3)$ | $O(n^4)$ | $O(n^4)$ | |
| **alignment** | time | $O(n^2 m^2)$ | $O(mn^4)$ | $O(mn^4)$ | $O(mn^5)$ | $O(mn^5)$ | $O(mn^6)$ | NP-complete |
| | space | $O(nm)$ | $O(mn^2)$ | $O(mn^3)$ | $O(mn^3)$ | $O(mn^4)$ | $O(mn^4)$ | |

**new contribution**

- implementation: PKalign
- accepted for RECOMB 09

# Literature (1/2)

**[Akutsu2000]** Akutsu, T. Dynamic programming algorithms for RNA secondary structure prediction with pseudoknots. *Discrete Applied Mathematics,* **2000***, 104*, 45-62

**[Bauer2007]** Bauer, M.; Klau, G. W. & Reinert, K. Accurate multiple sequence-structure alignment of RNA sequences using combinatorial optimization *BMC Bioinformatics,* 2007*, 8*, 271

**[Dirks2003]** Dirks, R. M. & Pierce, N. A. A partition function algorithm for nucleic acid secondary structure including pseudoknots. *J Comput Chem,* 2003*, 24*, 1664-77

**[Evans1999]** Evans, P. A. Finding Common Subsequences with Arcs and Pseudoknots. *CPM '99: Proceedings of the 10th Annual Symposium on Combinatorial Pattern Matching, Springer-Verlag,* 1999, 270-280

**[Evans2006]** Evans, P. A. Finding Common RNA Pseudoknot Structures in Polynomial Time. *Combinatorial Pattern Matching (CPM 2006), Springer Berlin / Heidelberg,* 2006*, 4009/2006*, 223-232

**[Jiang2002]** Tao Jiang, Guohui Lin, Bin Ma, and Kaizhong Zhang. A general edit distance between RNA structures. J. Comput. Biol., 9(2):371–88, 2002.

**[Lyngso2000]** Lyngso, R. B. & Pedersen, C. N. S. Pseudoknots in RNA Secondary Structures. *RECOMB 00, ACM Press,* 2000

**[Lenhof1998]** Lenhof, H. P.; Reinert, K. & Vingron, M. A Polyhedral Approach to RNA Sequence Structure Alignment 1998

**[Missal2005]** Missal, K.; Rose, D. & Stadler, P. F. Non-coding RNAs in Ciona intestinalis. *Bioinformatics,* 2005*, 21 Suppl 2*, ii77-ii78

# Literature (2/2)

**[Möhl2008]** Möhl, M.; Will, S. & Backofen, R. Fixed Parameter Tractable Alignment of RNA Structures. *Proceedings of the 19th Annual Symposium on Combinatorial Pattern Matching (CPM 2008),* 2008

**[Needleman1970]** Needleman, S. B. & Wunsch, C. D. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J Mol Biol,* 1970*, 48*, 443-453

**[Notredame2000]** Notredame, C.; Higgins, D. G. & Heringa, J. T-Coffee: A novel method for fast and accurate multiple sequence alignment. J Mol Biol, 2000, 302, 205-17

**[Reeder2004]** Reeder, J. & Giegerich, R. Design, implementation and evaluation of a practical pseudoknot folding algorithm based on thermodynamics. *BMC Bioinformatics,* 2004*, 5*, 104

**[Rivas1999]** Rivas, E. & Eddy, S. R. A dynamic programming algorithm for RNA structure prediction including pseudoknots. *J Mol Biol,* 1999*, 285*, 2053-2068

**[Sankoff1985]** Sankoff, D. Simultaneous solution of the RNA folding, alignment and protosequence problems. *SIAM J. Appl. Math.,* 1985*, 45*, 810-825

**[Siebert2007]** Siebert, S. & Backofen, R. Methods for multiple alignment and consensus structure prediction of RNAs implemented in MARNA *Methods Mol Biol,* **2007***, 395*, 489-502

**[Uemura1999]** Uemura, Y.; Hasegawa, A.; Kobayashi, S. & Yokomori, T. Tree adjoining grammars for RNA structure prediction. T*heoretical Computer Science,* 1999*, 210*, 277 - 303

**[Zuker 1981]** Zuker, M. & Stiegler, P. Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information. *Nucleic Acids Res,* 1981*, 9*, 133-148