

# Seminar presentation

## RNAalifold: improved consensus structure prediction for RNA alignments

Stephan H Bernhart et al., 2008, BMC Bioinformatics

# Introduction

---

- About 90% of mammalian genome is transcribed
- < 2% is protein-coding
- Most of the rest is probably functional RNA (tRNA, microRNA, snRNA etc.)
  - Huge amount of RNA molecules to analyse

# Introduction

---

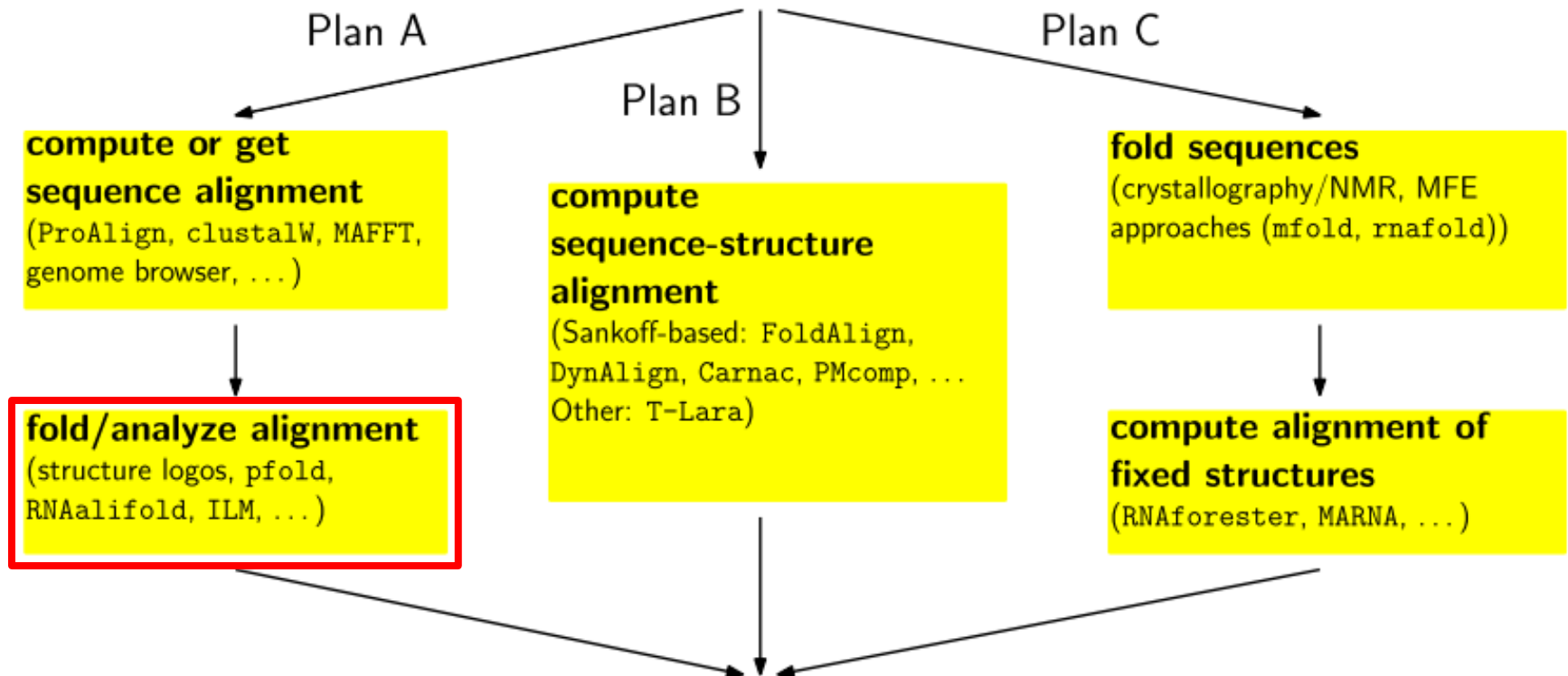
- Most functional RNAs have characteristic secondary structures
  - Highly conserved during evolution
- „Similiar structure“ possibly leads to „similar function“
- Structure is higher conserved than underlying sequences
  - RNA sequences vary between species but structure usually does not

**Task:** Given a sequence, find the optimal secondary structures or  
Given a MSA, find the consensus structure.

# Introduction

```
UGCAGCGACGGAAACGCUGCUAGCUUUGCGGCUAAGACUCUCGA
CGUGCCGAAAUGGCCGGGGCUCCACCGAGGAUGAUGC
ACGAUGAUGAUCGAUCGAUCGGACGUAGCUGACUAGCUGACU
```

Homologous RNA sequences



```
UGCAGCGACGGAAACG--CUGC----UAGCUUUGCGGCUAAGACUCUCGA
.((((((.....))--)))----((((.....)))).....
ACGAUGAU-GAUCGAUCGAUCGGACGUAGCUGACUAGCUGA---CU----
```

source: <http://www.mi.fu-berlin.de/wiki/pub/ABI/SS14Lecture11Materials/script.pdf>, SeqAna RNA script

# Introduction

RNAalifold is a DP algorithm, that uses phylogenetic information (sequence covariance) and thermodynamic stability of molecules to predict a consensus structure of a MSA.

## Idea behind sequence covariance

*example:*

seq1	G	<b>C</b>	C	U	U	C	C	C	<b>G</b>	C
seq2	G	<b>A</b>	C	U	U	C	C	C	<b>U</b>	C
seq3	G	<b>G</b>	C	U	U	C	C	C	<b>C</b>	C

Bold printed basepairs are covarying (compensatory mutations).

→ base pair most likely present in consensus structure

# Methods

## The old RNAalifold

- 4 matrices corresponding to different structural components:
  - F (unconstrained structures)
  - C (constrained structures)
  - M (multi-loop structures)
  - M<sup>1</sup>(multi-loop with one branch)
- They hold for every sub-sequence from i to j the optimal folds
- $\gamma(i,j)$  is the covariance score

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

$$C_{i,j} = \beta\gamma(i,j) + \min \left\{ \begin{array}{l} \sum_{\alpha \in \mathbb{A}} \mathfrak{H}(i,j,\alpha) \\ \min_{i < k < l < j} \left( \sum_{\alpha \in \mathbb{A}} \mathfrak{J}(ij,kl,\alpha) + C_{k,l} \right) \\ \min_{i < k < j} \left( M_{i,k} + M_{k+1,j}^1 + \mathbf{a} \right) \end{array} \right.$$

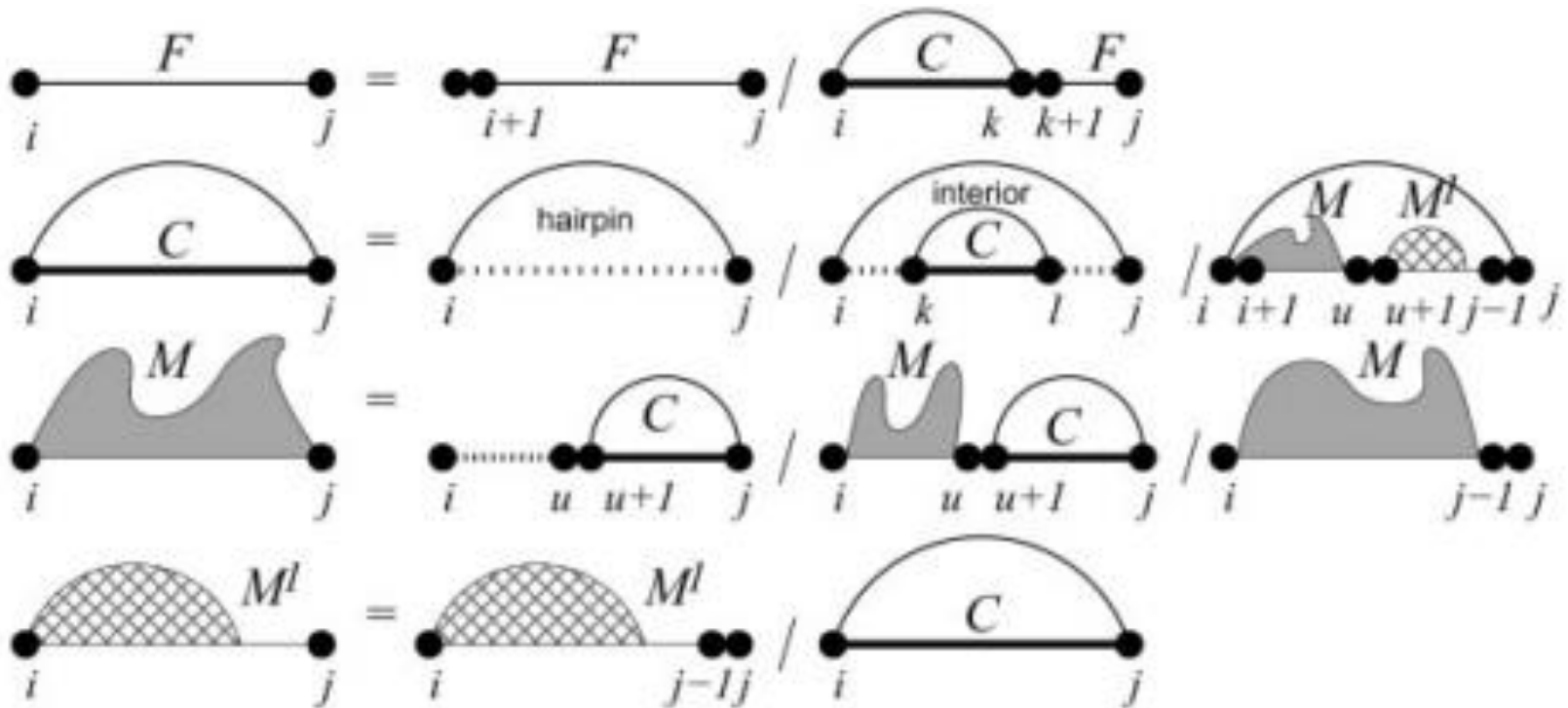
$$M_{i,j} = \min \left\{ \begin{array}{l} M_{i+1,j} + \mathbf{c} \\ \min_{i < k < j} C_{i,k} + M_{k+1,j} + \mathbf{b} \\ M_{i,j}^1 \end{array} \right.$$

$$M_{i,j}^1 = \min \left( M_{i,j-1}^1 + \mathbf{c}, C_{i,k} \right)$$

source: Bernhart et al., 2008, BMC Bioinformatics

# Methods

## „Visualization“ of the recursions



source: Hofacker/Stadler, 2006

# Methods

## Improvement of the old RNAalifold

**old covariance score:**

$$\gamma'(i, j) = \frac{1}{2} \sum_{\substack{\alpha, \beta \in \mathbb{A} \\ \alpha \neq \beta}} \begin{cases} h(\alpha_i, \beta_i) + h(\alpha_j, \beta_j) & \text{if } (\alpha_i, \alpha_j) \in \mathcal{B} \\ & \wedge (\beta_i, \beta_j) \in \mathcal{B} \\ 0 & \text{otherwise} \end{cases}$$

$$\mathcal{B} = \{AU, UA, CG, GC, GU, UG\}$$

- Based on hamming distance  $h(a, b)$
- **Problem:** No quantitative argument, since  $h(a, b)$  is either 1 or 0. Some mutations might be more frequent than other...



# Methods

## Improvement of the old RNAalifold

- **Solution:** Introducing a scoring matrix (RIBOSUM)
- Scores derived from frequencies of aligned and basepaired nucleotides (log odds)
- example: RIBOSUM85-60

AA	-2.49																
AC	-7.04	-2.11															
AG	-8.24	-8.89	-0.80														
AU	-4.32	-2.04	-5.13	<b>4.49</b>													
CA	-8.84	-9.37	-10.41	-5.56	-5.13												
CC	-14.37	-9.08	-14.53	-6.71	-10.45	-3.59											
CG	-4.68	-5.86	-4.57	<b>1.67</b>	-3.57	-5.71	<b>5.36</b>										
CU	-12.64	-10.45	-10.14	-5.17	-8.49	-5.77	-4.96	-2.28									
GA	-6.86	-9.73	-8.61	-5.33	-7.98	-12.43	-6.00	-7.71	-1.05								
GC	-5.03	-3.81	-5.77	<b>2.70</b>	-5.95	-3.70	<b>2.11</b>	-5.84	-4.88	<b>5.62</b>							
GG	-8.39	-11.05	-5.38	-5.61	-11.36	-12.58	-4.66	-13.69	-8.67	-4.13	-1.98						
GU	-5.84	-4.72	-6.60	<b>0.59</b>	-7.93	-7.88	-0.27	-5.61	-6.10	<b>1.21</b>	-5.77	<b>3.47</b>					
UA	-4.01	-5.33	-5.43	<b>1.61</b>	-2.42	-6.88	<b>2.75</b>	-4.72	-5.85	<b>1.60</b>	-5.75	-0.57	<b>4.97</b>				
UC	-11.32	-8.67	-8.87	-4.81	-7.08	-7.40	-4.91	-3.83	-6.63	-4.49	-12.01	-5.30	-2.98	-3.21			
UG	-6.16	-6.93	-5.94	-0.51	-5.63	-8.41	<b>1.32</b>	-7.36	-7.55	-0.08	-4.27	-2.09	<b>1.14</b>	-4.76	<b>3.36</b>		
UU	-9.05	-7.83	-11.07	-2.98	-8.39	-5.41	-3.67	-5.21	-11.54	-3.90	-10.79	-4.45	-3.39	-5.97	-4.28	-0.02	
	AA	AC	AG	AU	CA	CC	CG	CU	GA	GC	GG	GU	UA	UC	UG	UU	

source: Klein/Eddy, 2003, BMC Bioinformatics

# Methods

## Improvement of the old RNAalifold

- **new covariance score:** 
$$\gamma'(i, j) = \frac{1}{2} \sum_{\substack{\alpha, \beta \in A \\ \alpha \neq \beta}} xR(\alpha_i \alpha_j; \beta_i \beta_j)$$

- Replacement of the hamming distance with the RIBOSUM scores
- *Complete covariance score used in recursions:*

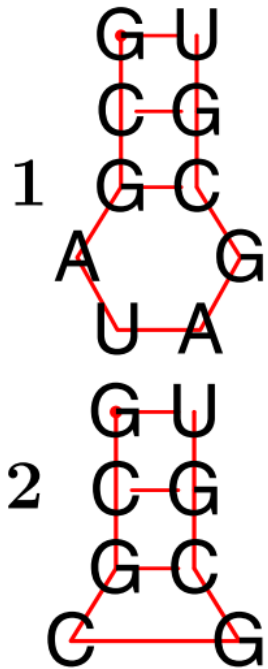
$$\beta \cdot \gamma(i, j) = \beta \cdot \left[ \gamma'(i, j) + \delta \sum_{\alpha \in A} \begin{cases} 0, & \text{if } (\alpha_i \alpha_j) \in B \\ 0.25, & \text{if } \alpha_i \text{ and } \alpha_j \text{ are gaps} \\ 1, & \text{otherwise} \end{cases} \right]$$

- Parameter  $\beta$  – influence of covariance on folding
- Parameter  $\delta$  – impact of non-standard base pairs

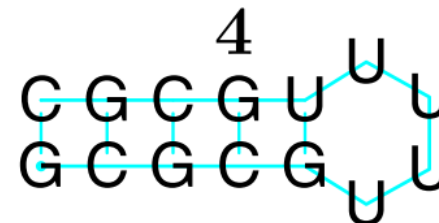
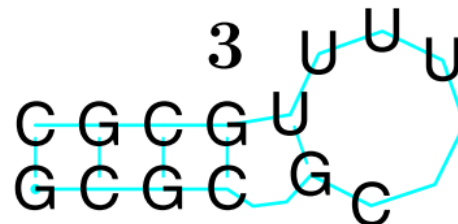
# Methods

## Improvement of the old RNAalifold

Removal of gaps in interior loops before calculating energies  
 → less energetic unfavorable loops



		*****	*	*****		
sequence_2	AGCGUUCUUGC	CGC	--	GUGUUUUUGC	CGCUUGC	30
sequence_3	AGCGUUCUUGC	CGC	--	GU--UUUUGC	CGCUUGC	28
sequence_1	AGCGUUCUUGC	GAUAG		CGUUUUUGC	CGCUUGC	32
old		((((((.....(((.....))).....)))))).....				-5.95
new		((((((.....(((.....))).....)))))).....				-5.83



# Results

---

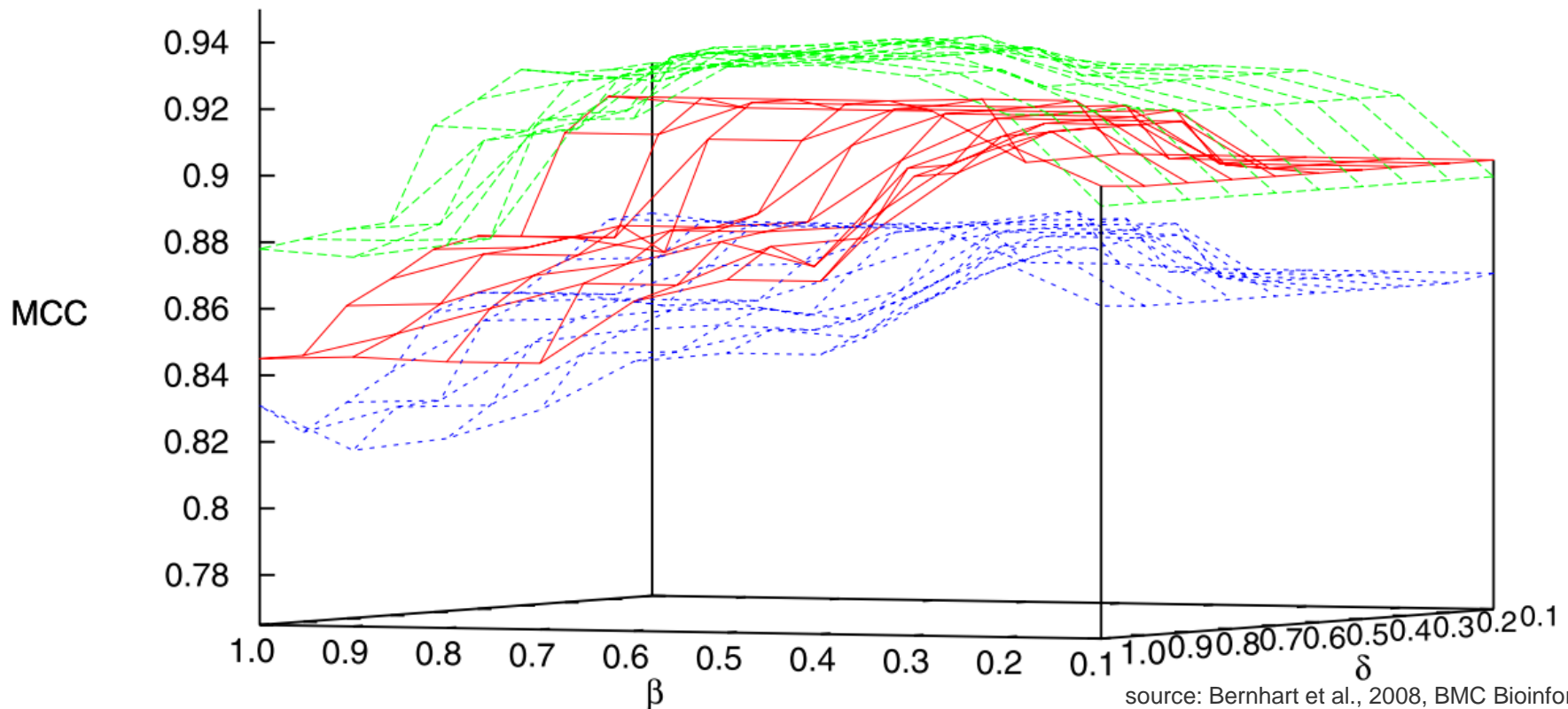
## Measuring accuracy

- CMfinder-SARSE dataset for reference structures
- Mathews correlation coefficient (MCC) for accuracy measuring
  - Basically a value to compare reference structure with predicted structure
    - $MCC = 1$ , perfect prediction
    - $MCC = 0$ , not better than random
    - $MCC = -1$ , total disagreement

# Results

## Influence of $\beta$ and $\delta$ on accuracy

RNAalifold NEW ———  
 with RIBOSUM - - -  
 RNAalifold OLD ·····



source: Bernhart et al., 2008, BMC Bioinformatics

# Results

---

## Measuring accuracy

- New RNAalifold RIBOSUM better than other two (except for  $\beta = 1$ )
- Determination of optimal values for  $\beta$  and  $\delta$  (0.6 and 0.5, respectively)
- Old RNAalifold emphasized covariance too much

# Results

## Comparison to other Folding algorithms

RNA	#seq	MPI	RIBOSUM	RNAalifold	Pfold	KNetFold	McC_mea
Antizyme_FSE	13	87	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
ctRNA_pGAI	15	72	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	0.976	<b>1.000</b>
Entero_5_CRE	160	84	<b>1.000</b>	0.848	0.478	<b>1.000</b>	0.942
Entero_CRE	56	81	<b>1.000</b>	0.736	<b>1.000</b>	0.953	0.953
GcvB	17	64	<b>0.939</b>	0.799	0.889	<b>0.939</b>	0.921
glmS	11	60	<b>0.986</b>	0.972	0.972	0.809	0.837
HACA_sno_Snake	22	90	0.871	0.407	0.414	<b>0.915</b>	0.884
HCV_SLIV	110	89	<b>1.000</b>	0.922	<b>1.000</b>	<b>1.000</b>	0.961
HDV_ribozyme	15	95	<b>0.953</b>	-0.015	0.590	0.460	0.460
HepC_CRE	52	87	<b>1.000</b>	0.962	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
Histone3	64	78	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
Hsp90_CRE	4	98	0.855	0.855	0.413	0.867	<b>0.874</b>
IBV_D-RNA	10	96	<b>1.000</b>	0.928	0.928	<b>1.000</b>	<b>1.000</b>
Intron_gpII	114	54	<b>1.000</b>	0.779	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
[.....]							
SNORD64	3	94	<b>1.000</b>	0.539	0.539	0.661	-0.014
SNORD86	6	82	<b>0.641</b>	-0.012	-0.007	0.511	0.000
snoU83B	4	87	<b>0.927</b>	<b>0.927</b>	0.846	0.895	<b>0.927</b>
TCV_H5	3	97	<b>1.000</b>	<b>1.000</b>	0.685	<b>1.000</b>	<b>1.000</b>
TCV_Pr	4	95	<b>1.000</b>	<b>1.000</b>	0.688	<b>1.000</b>	<b>1.000</b>
Tymo_tRNA-like	28	64	<b>1.000</b>	0.916	<b>1.000</b>	0.973	<b>1.000</b>
ykoK	36	61	0.856	0.756	<b>0.906</b>	0.841	0.794
mean			<b>0.937</b>	0.831	0.765	0.866	0.837

# Discussion

## Comparison with other folding algorithms (RNA STRAND – Rfam set)

RNA	comment	RIBOSUM	RNAalifold	Pfold	KNetFold	McC_mea
7SK		<b>0.507</b>	0.456	0.292	0.429	0.306
bicoid_3		<b>0.949</b>	0.840	n.a.	0.829	0.927
Corona_pk3	Pk	0.579	0.646	0.674	0.678	<b>0.705</b>
CPEB3_ribozyme	Pk	<b>0.756</b>	<b>0.756</b>	0.663	<b>0.756</b>	0.612
Gammaretro_CES		<b>0.983</b>	0.948	<b>0.983</b>	0.935	<b>0.983</b>
Hammerhead_I		<b>1.000</b>	0.474	0.621	0.831	0.614
Hammerhead_3		<b>1.000</b>	0.960	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
HDV_ribozyme	Pk	0.709	-0.018	<b>0.784</b>	0.388	0.396
IRES_c-myc		-0.004	0.079	0.286	-0.002	<b>0.350</b>
[.....]						
Telomerase-vert	pk	<b>0.918</b>	0.751	n.a.	n.a.	0.820
Vimentin3		0.741	-0.016	0.184	<b>0.771</b>	0.629
Y		<b>1.000</b>	<b>1.000</b>	0.925	<b>1.000</b>	<b>1.000</b>
mean		0.759	0.651			0.703
mean	knetfold	0.750	0.645		0.680	0.696
mean	pfold	0.756	0.635	0.693	0.682	0.673

source: Bernhart et al., 2008, BMC Bioinformatics



# Results

---

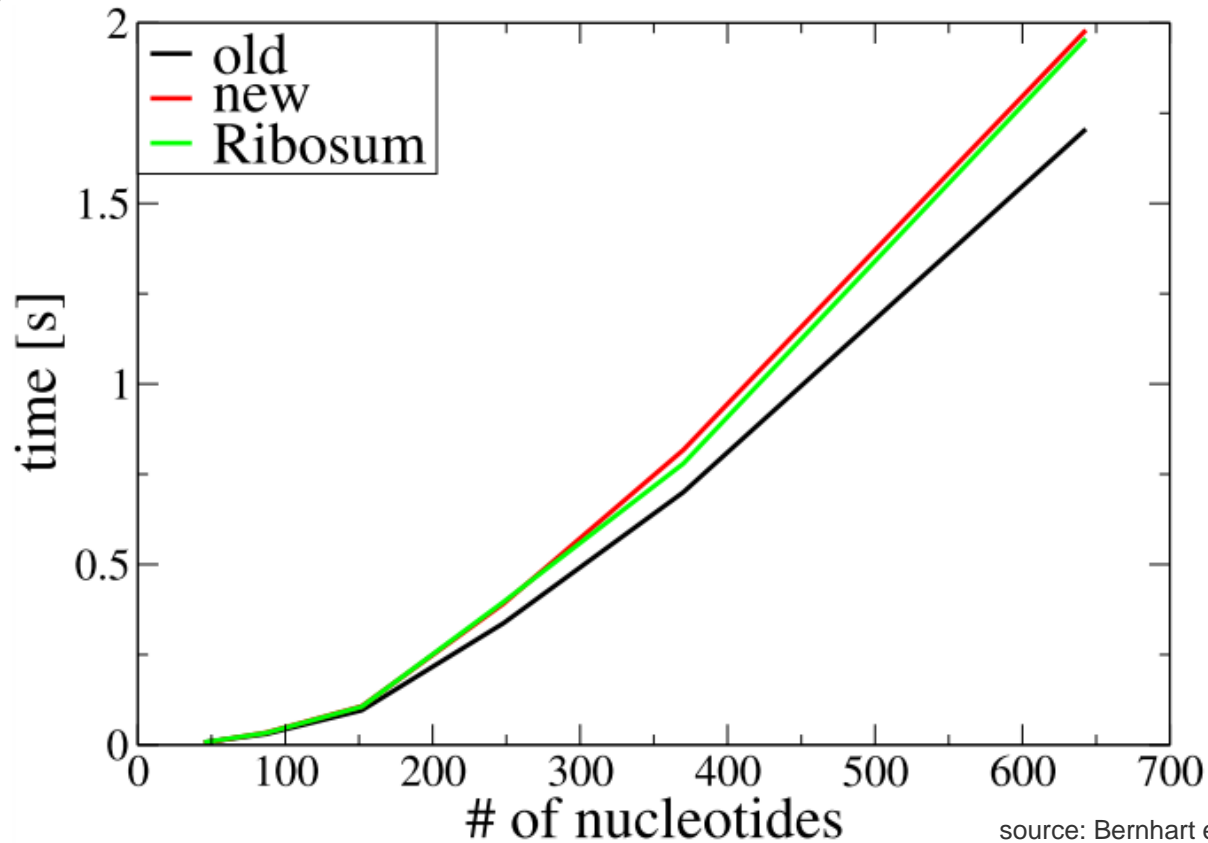
## Computational performance

- Runtime:  $O(Nn^3)$ ,  $N$  sequences of length  $n$
- Memory:  $O(n^2)$
- Tests run on an Intel Xeon 2.8 GHz
- Only runtime was tested

# Results

## Measuring runtime

- $N = 4$
- $n$  variable

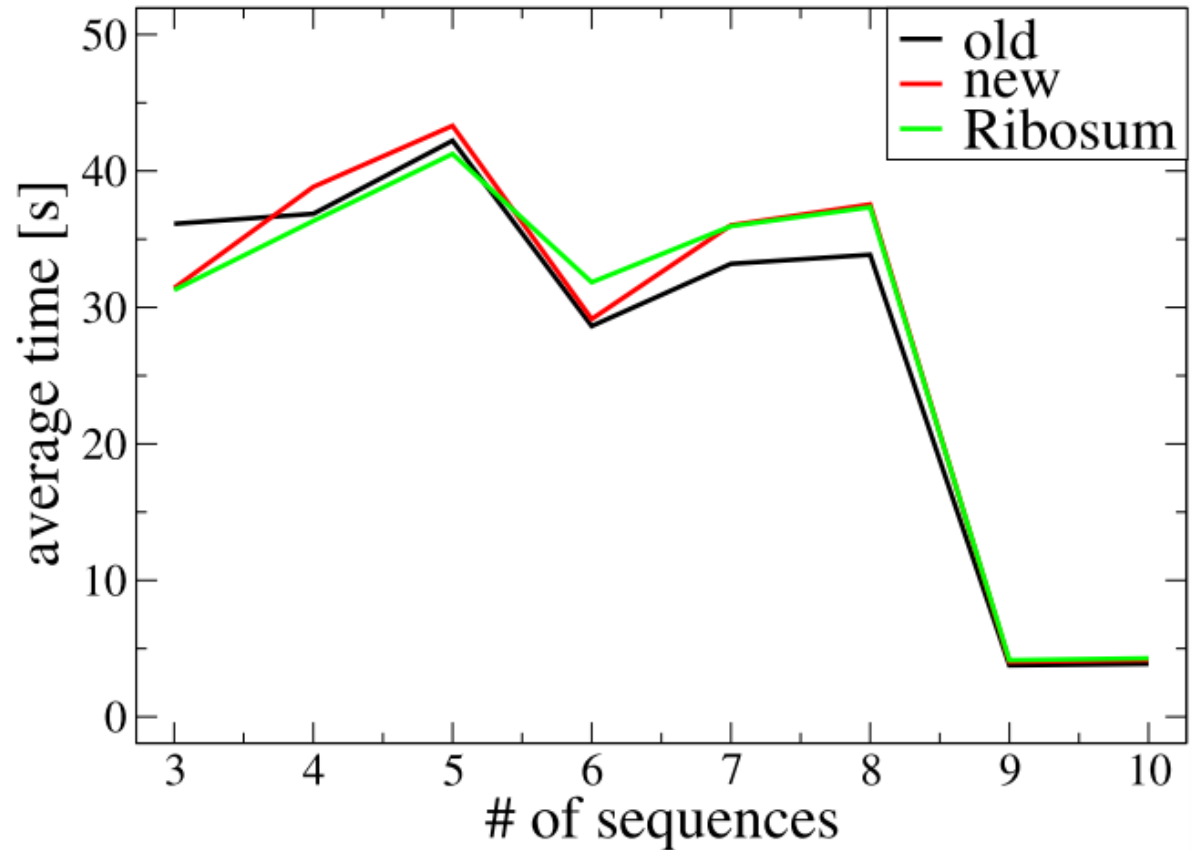


source: Bernhart et al., 2008, BMC Bioinformatics

# Results

## Measuring runtime

- $n = 1716$
- $N$  variable
- All sequences from same alignment.



source: Bernhart et al., 2008, BMC Bioinformatics

Dropping basepairs that aren't supported by sequences (cutoff).

→ The more sequences available the less basepairs have to be considered.

# Fazit

---

- Changes to gaps and covariance led to improved accuracy
- Slightly more computational effort
- The new RNAalifold is on par or better than comparative prediction software (at least on tested data sets)
- No pseudoknots (DP algorithms cannot detect pk)

**Thanks for your attention!**