

SPARSE: Quadratic Time SA&F of RNAs without Sequence-Based Heuristics

Beijing, RECOMB 2013

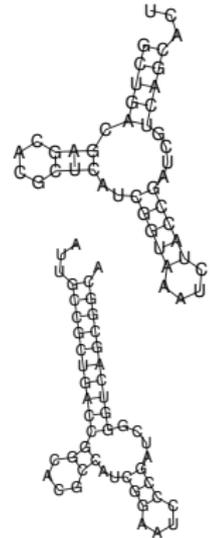
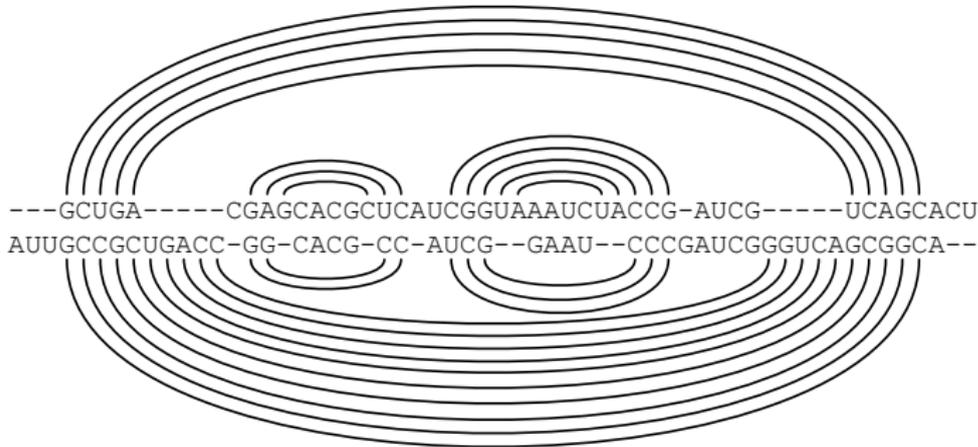
Sebastian Will, Christina Schmiedl, Milad Miladi,
Mathias Möhl, Rolf Backofen

University of Leipzig
University of Freiburg

Simultaneous Alignment and Folding [Sankoff]

Given: A = GCUGACGAGCACGCUCAUCGGUAAAUCUACCGAUCGUCAGCACU
& B = AUUGCCGCUGACC-GG-CACG-CC-AUCG--GAU--CCCGAUCGGGUCAGCGGCA--

Find:



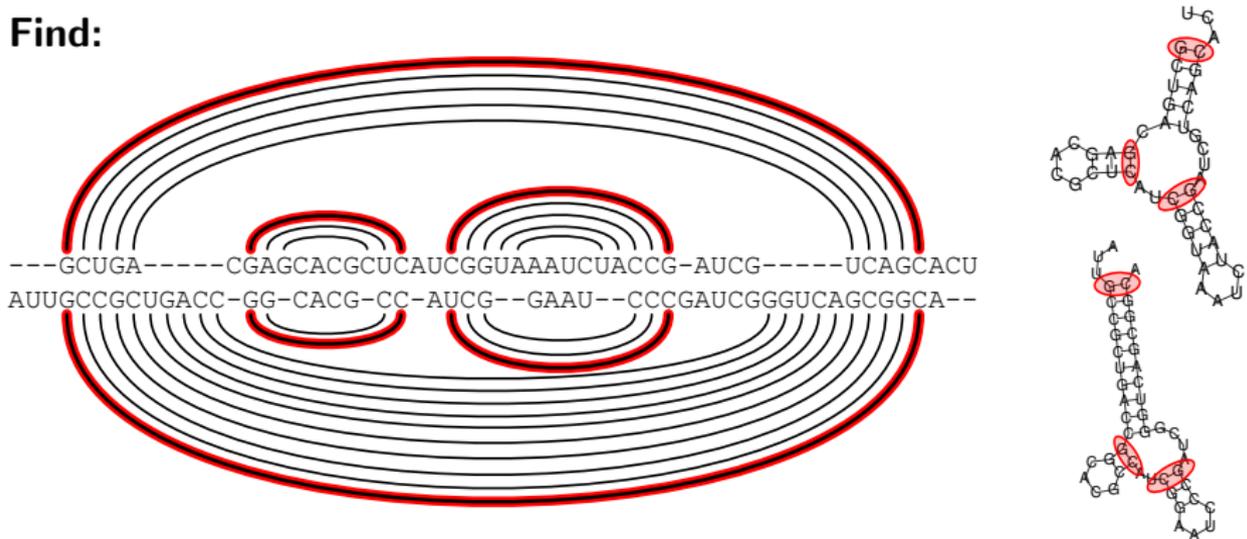
sequence similarity + energy A + energy B \rightarrow opt

where alignment, structure A, & structure B are **compatible**

Simultaneous Alignment and Folding [Sankoff]

Given: A = GCUGACGAGCACGCUCAUCGGUAAAUCUACCGAUCGUCAGCACU
& B = AUUGCCGCUGACC-GG-CACG-CC-AUCG--GAAU--CCCGAUCGGGUCAGCGGCA

Find:



sequence similarity + energy A + energy B \rightarrow opt

where alignment, structure A, & structure B are **compatible**

Sankoff's Algorithm

Dynamic Programming

RNA Energy Minimization [Zuker]

×

Sequence Alignment

$O(n^6)$ = “extreme computational cost”

Sankoff's Algorithm

Dynamic Programming

RNA Energy Minimization [Zuker]

×

Sequence Alignment

$O(n^6)$ = “extreme computational cost”

Sankoff's Algorithm

Dynamic Programming

RNA Energy Minimization [Zuker]

×

Sequence Alignment

$O(n^6)$ = “extreme computational cost”

Sankoff-style Approaches

HEAVY

Dynalign
FoldAlign

- Sankoff implementations
- heavyweight energy model
- sequence-based heuristics

LIGHT

PMcomp

- lightweight energy model
- base pair probabilities

LocARNA

- + sparsifies structure space (ensemble-based)
- improves time and space

RAF

- + sparsifies alignment space
- sequence-based heuristics

SPARSE

- strong sparsification w/o sequence-based heuristics

Sankoff-style Approaches

HEAVY

Dynalign
FoldAlign

- Sankoff implementations
- heavyweight energy model
- sequence-based heuristics

LIGHT

PMcomp

- lightweight energy model
- base pair probabilities

LocARNA

- + sparsifies structure space (ensemble-based)
- improves time and space

RAF

- + sparsifies alignment space
- sequence-based heuristics

SPARSE 

- strong sparsification w/o sequence-based heuristics

PMcomp's Trick – Lightweight SAF

Sankoff: **sequence similarity**
+ energies of A and B → **opt**

- **energies** composed of loop energies

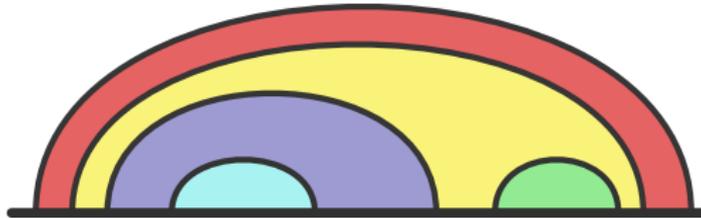


- Dynamic Programming
Base Pair Maximization [Nussinov] × Sequence Alignment
- **cheaper but same complexity**

PMcomp's Trick – Lightweight SAF

Sankoff: **sequence similarity**
+ energies of A and B → **opt**

- **energies** composed of loop energies

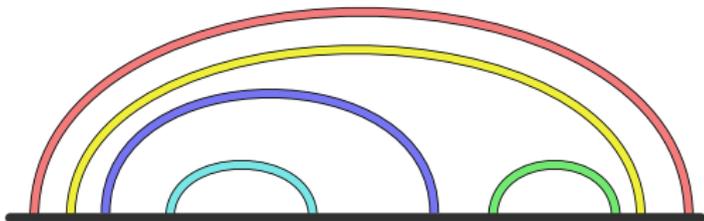


- Dynamic Programming
Base Pair Maximization [Nussinov] × Sequence Alignment
- **cheaper but same complexity**

PMcomp's Trick – Lightweight SAF

PMcomp: **sequence similarity**
+ **pseudo-energies of A and B** → **opt**

- **pseudo-energies** composed of “base pair energies”

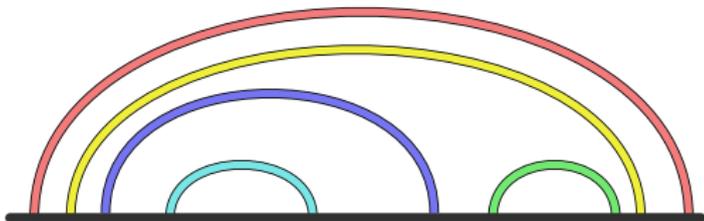


- Dynamic Programming
Base Pair Maximization [Nussinov] × Sequence Alignment
- **cheaper but same complexity**

PMcomp's Trick – Lightweight SAF

PMcomp: **sequence similarity**
+ **pseudo-energies of A and B** → **opt**

- **pseudo-energies** composed of “base pair energies”

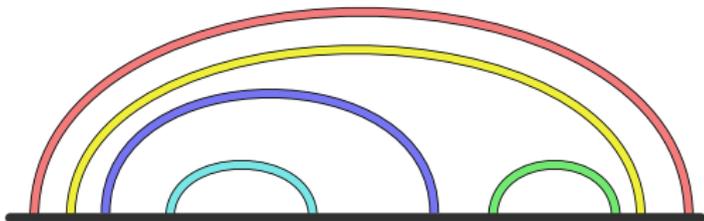


- Dynamic Programming
Base Pair Maximization [Nussinov] × Sequence Alignment
- **cheaper but same complexity**

PMcomp's Trick – Lightweight SAF

PMcomp: **sequence similarity**
+ **pseudo-energies of A and B** → **opt**

- **pseudo-energies** composed of “base pair energies”



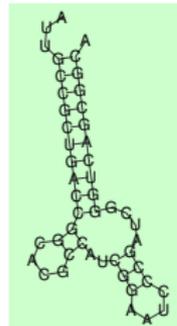
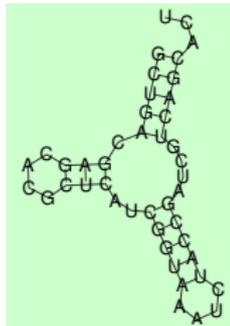
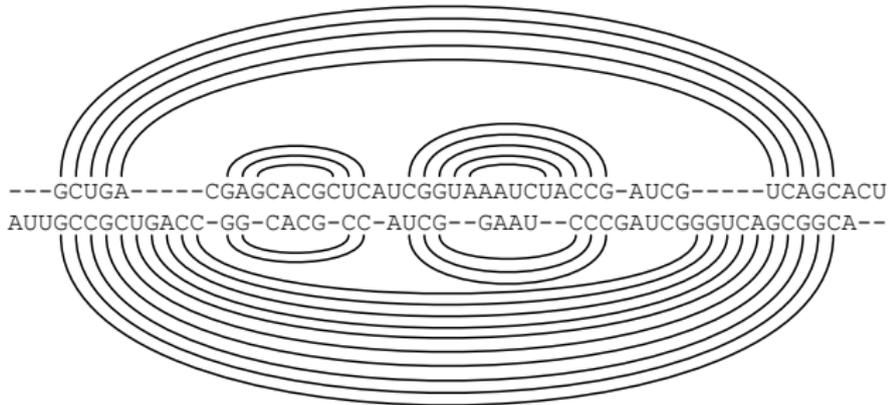
- Dynamic Programming
Base Pair Maximization [Nussinov] \times Sequence Alignment
- **cheaper but same complexity**

PMcomp – THE Lightweight Sankoff Algorithm?

compatibility

Sankoff: *same shape*

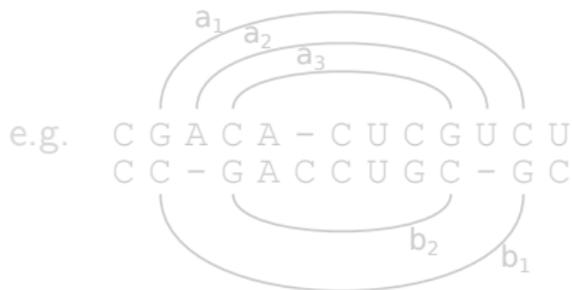
PMcomp: *all base pairs match*



PARSE — THE Lightweight Sankoff Algorithm

(PARSE = Prediction and Alignment of RNAs using Structure Ensembles)

- **lightweight** (PMcomp pseudo-energy)
& **complete** (Sankoff's compatibility)
- allows base pair **insertions and deletions**

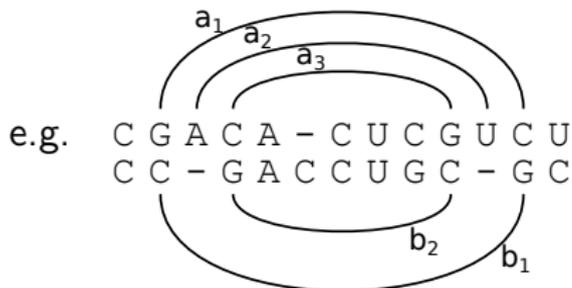


We need “complete” for strong sparsification, please be patient.

PARSE — THE Lightweight Sankoff Algorithm

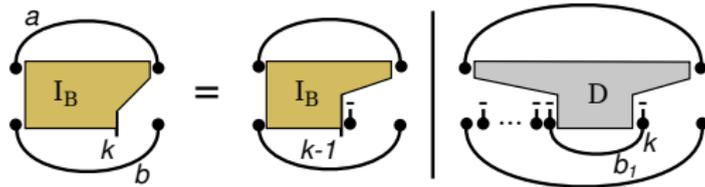
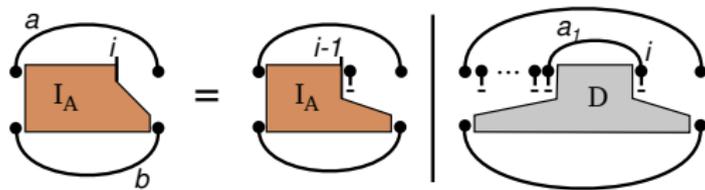
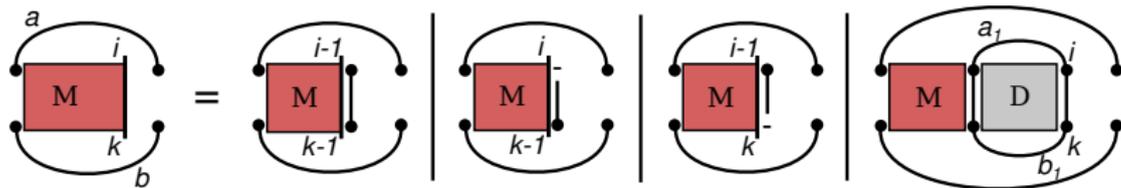
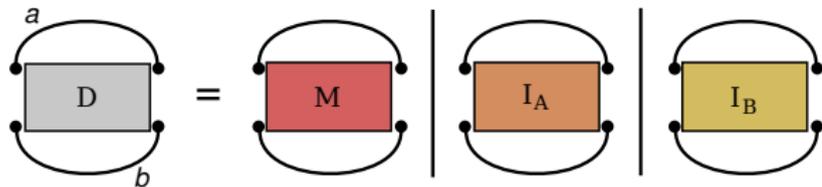
(PARSE = Prediction and Alignment of RNAs using Structure Ensembles)

- **lightweight** (PMcomp pseudo-energy)
& **complete** (Sankoff's compatibility)
- allows base pair **insertions and deletions**

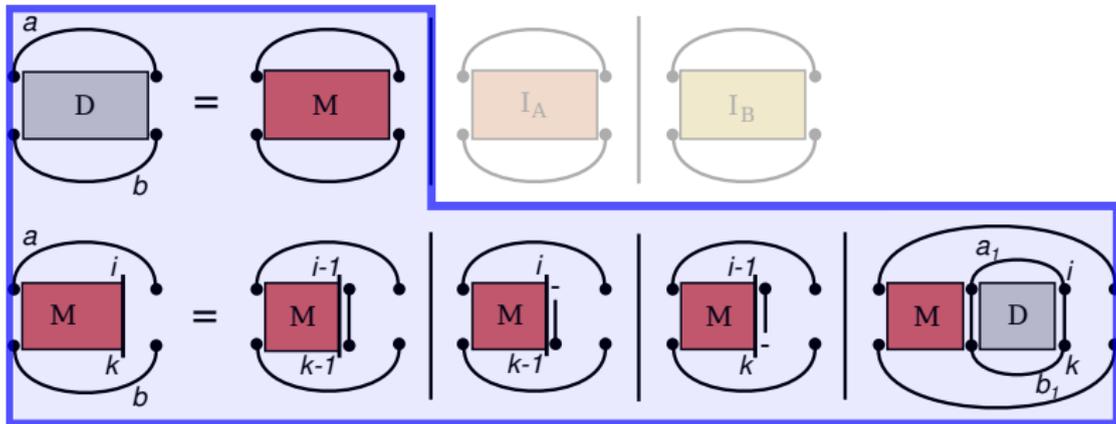


We need “complete” for strong sparsification, please be patient.

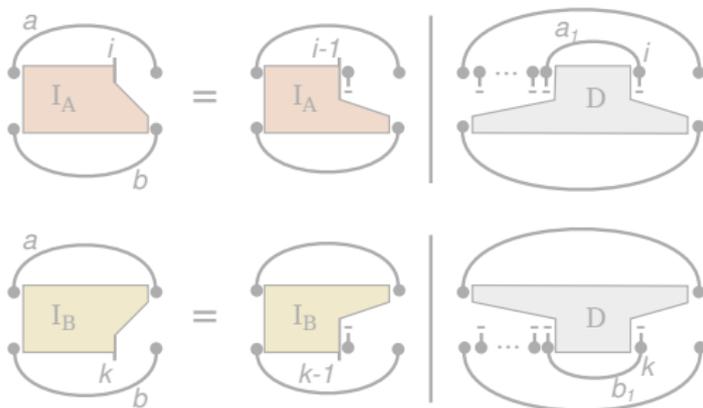
PARSE Algorithm



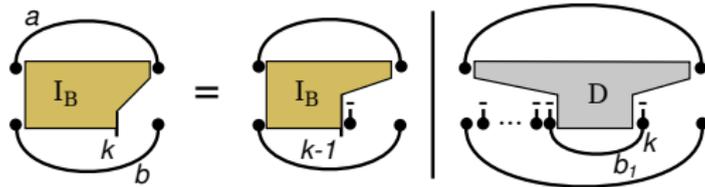
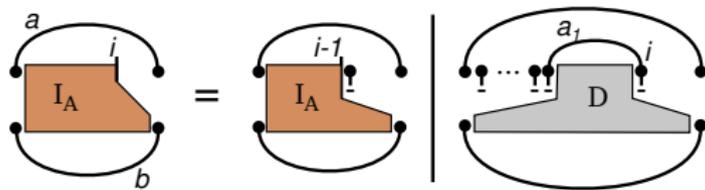
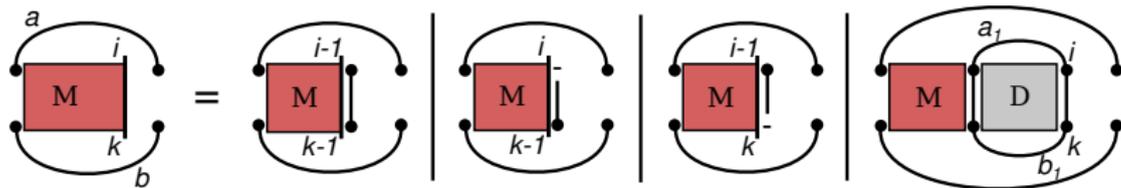
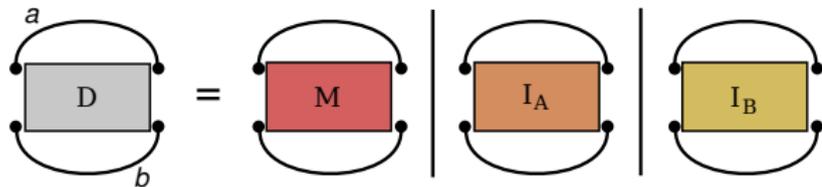
PARSE Algorithm



**PMcomp/
LocARNA-like
core**

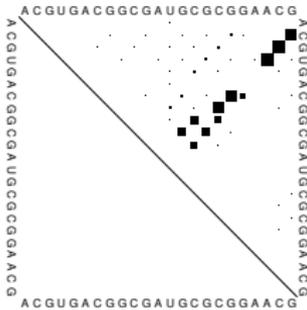


PARSE Algorithm



LocARNA's Trick: Ensemble-based Sparsification

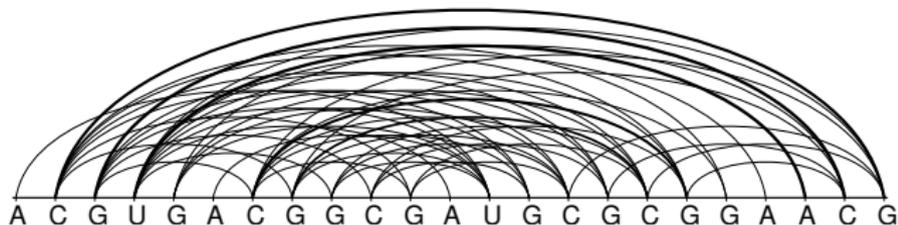
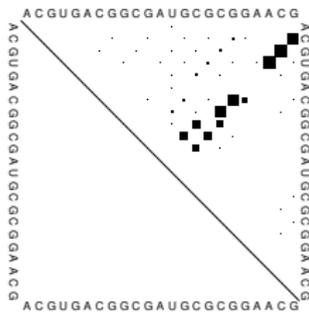
- Sparsify structure ensemble



- improves time and space; each by $O(n^2)$

LocARNA's Trick: Ensemble-based Sparsification

- Sparsify structure ensemble

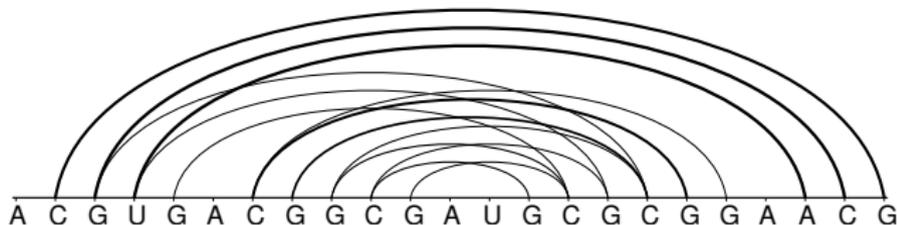
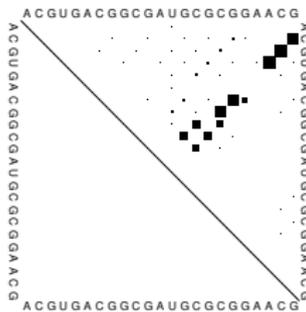


all base pairs

- improves time and space; each by $O(n^2)$

LocARNA's Trick: Ensemble-based Sparsification

- Sparsify structure ensemble

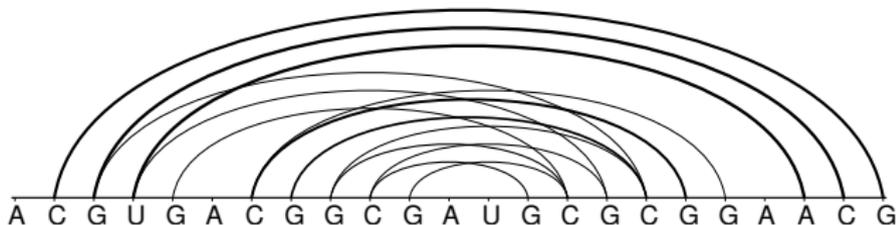
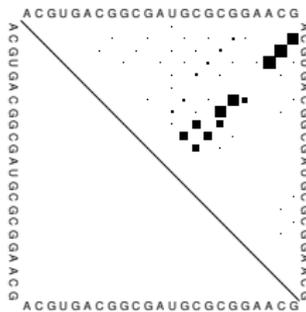


only probable base pairs

- improves time and space; each by $O(n^2)$

LocARNA's Trick: Ensemble-based Sparsification

- Sparsify structure ensemble



only probable base pairs

- improves time and space; each by $O(n^2)$

SPARSE: Novel Ensemble-based Sparsification*



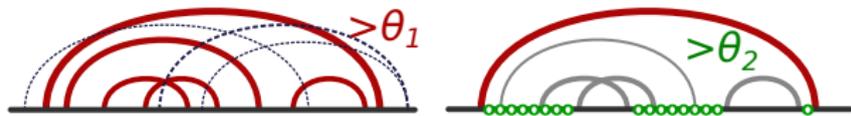
- only **base pairs** with probabilities $> \theta_1$
- only **bases** with unpaired probabilities in loops $> \theta_2$
- only **base pairs** with probabilities in loops $> \theta_3$

requires complete prediction (Sankoff/PARSE)

(*) confer LocARNA's "old" sparsification:

- match only base pairs with probabilities $> \theta_1$

SPARSE: Novel Ensemble-based Sparsification*



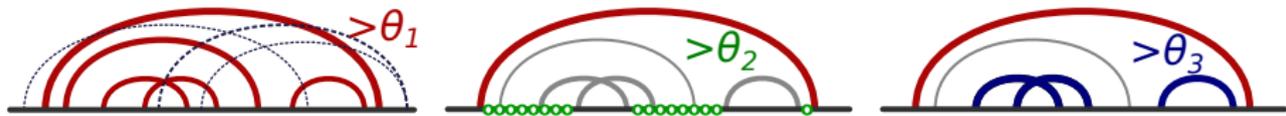
- only **base pairs** with probabilities $> \theta_1$
- only **bases** with unpaired probabilities in **loops** $> \theta_2$
- only **base pairs** with probabilities in **loops** $> \theta_3$

requires complete prediction (Sankoff/PARSE)

(*) confer LocARNA's "old" sparsification:

- match only base pairs with probabilities $> \theta_1$

SPARSE: Novel Ensemble-based Sparsification*



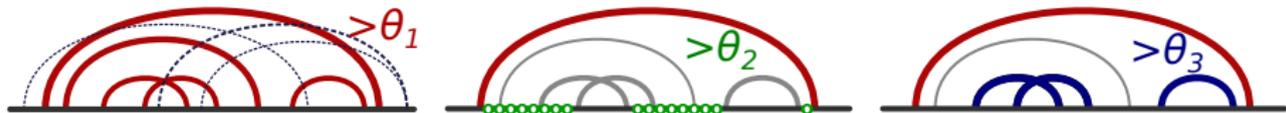
- only **base pairs** with probabilities $> \theta_1$
- only **bases** with unpaired probabilities in **loops** $> \theta_2$
- only **base pairs** with probabilities in **loops** $> \theta_3$

requires complete prediction (Sankoff/PARSE)

(*) confer LocARNA's "old" sparsification:

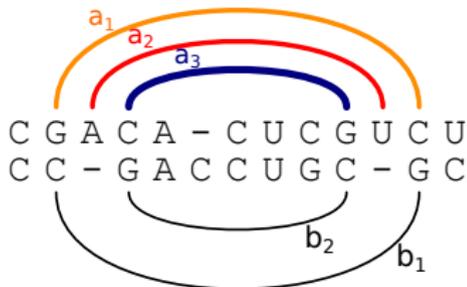
- match only base pairs with probabilities $> \theta_1$

SPARSE: Novel Ensemble-based Sparsification*



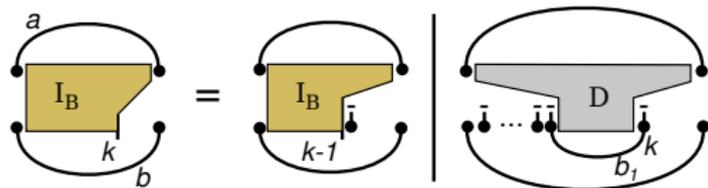
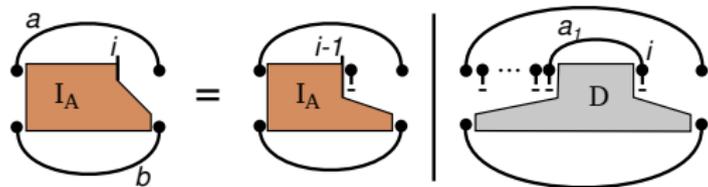
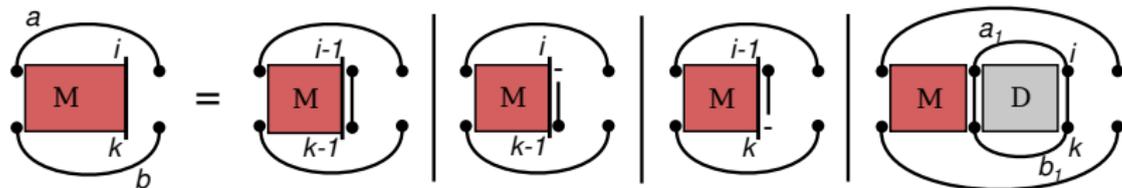
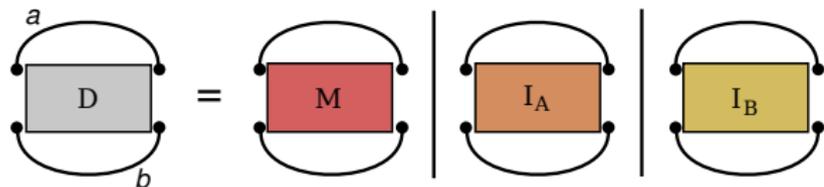
- only **base pairs** with probabilities $> \theta_1$
- only **bases** with unpaired probabilities in **loops** $> \theta_2$
- only **base pairs** with probabilities in **loops** $> \theta_3$

requires complete prediction (Sankoff/PARSE)

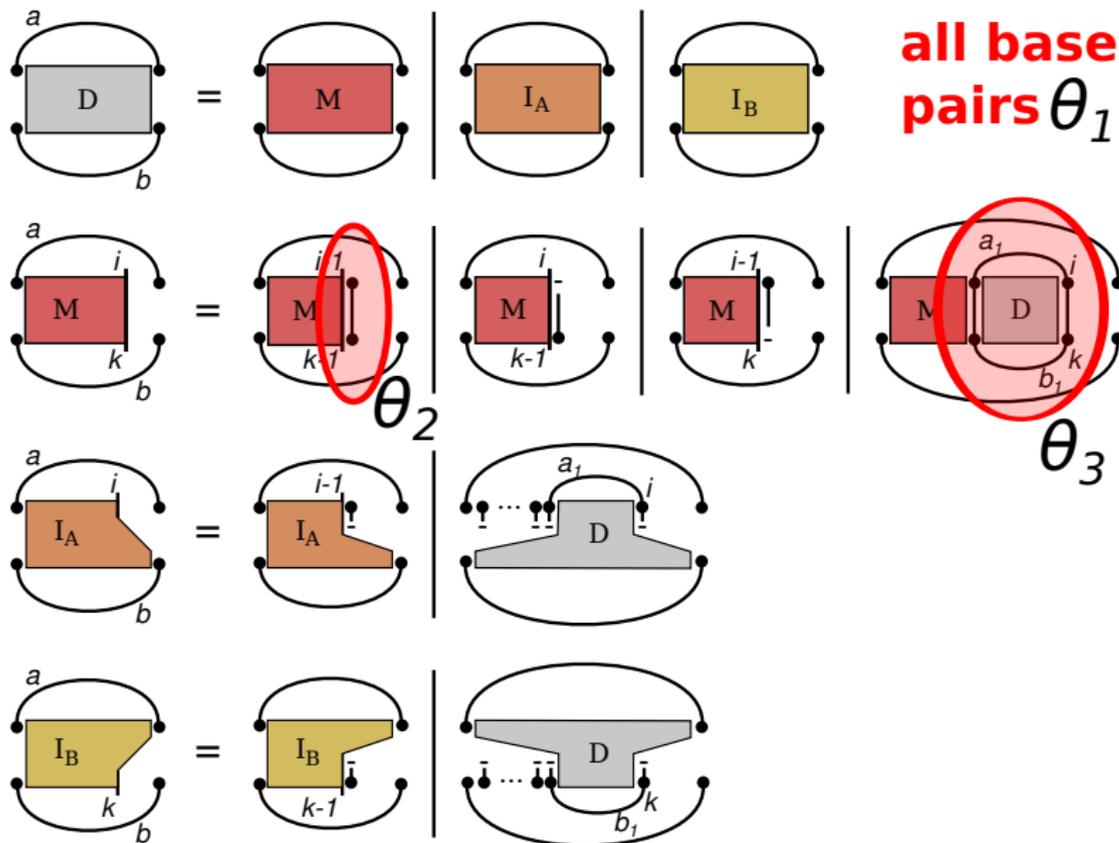


a_3 in loop a_2 ✓
 but a_3 in loop a_1 ✗
 a_2 ✗ \implies a_3 - b_2 ✗

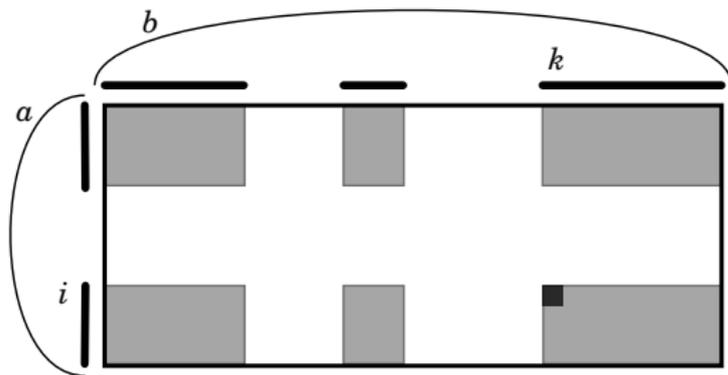
Thresholds in Recursions Cases



Thresholds in Recursions Cases

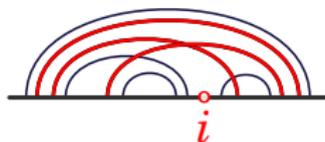


Quadratic Time



Q: How many matrices M^{ab} compute (i, k) ?

Count base pairs a where
 $\Pr^A[i \text{ in loop of } a] > \theta_2$

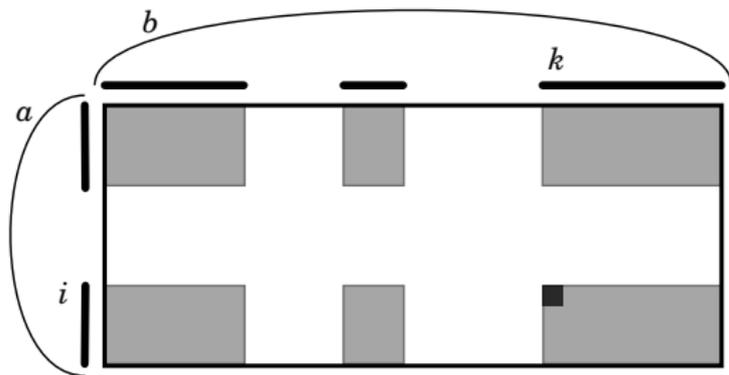


\Rightarrow less than $1/\theta_2$

A: each (i, k) in only constant number of matrices

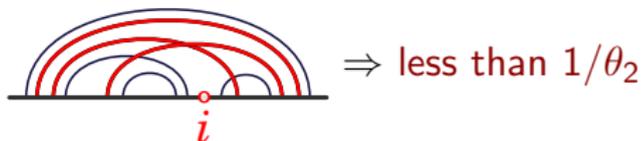


Quadratic Time



Q: How many matrices M^{ab} compute (i, k) ?

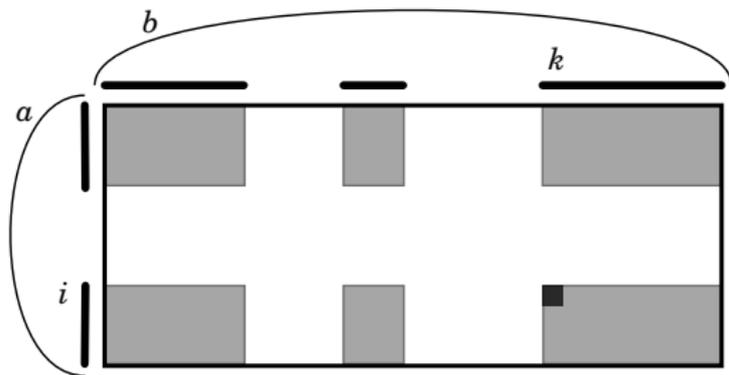
Count **base pairs** a where
 $\Pr^A[i \text{ in loop of } a] > \theta_2$



A: each (i, k) in only constant number of matrices

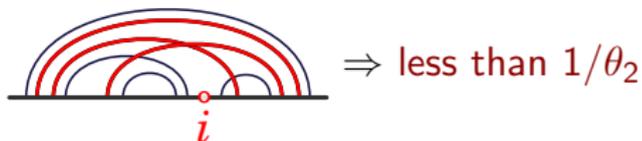


Quadratic Time



Q: How many matrices M^{ab} compute (i, k) ?

Count **base pairs** a where
 $\Pr^A[i \text{ in loop of } a] > \theta_2$



A: each (i, k) in only constant number of matrices

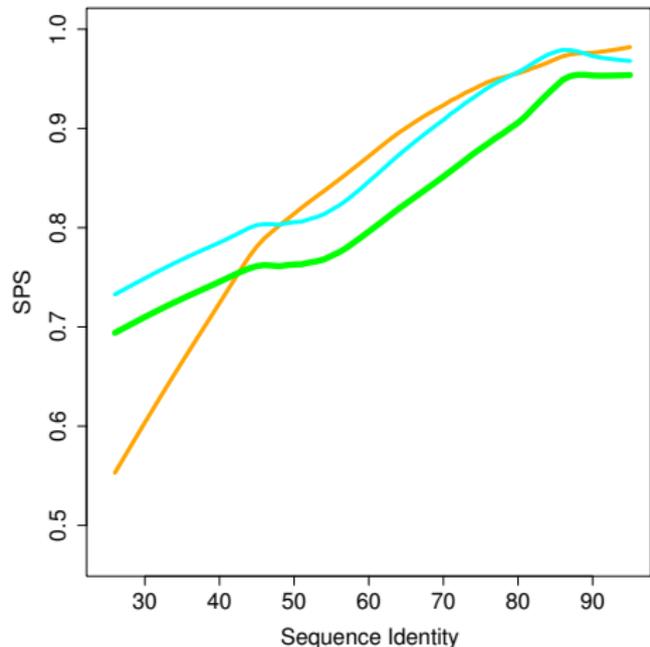


Run-times and Speedup

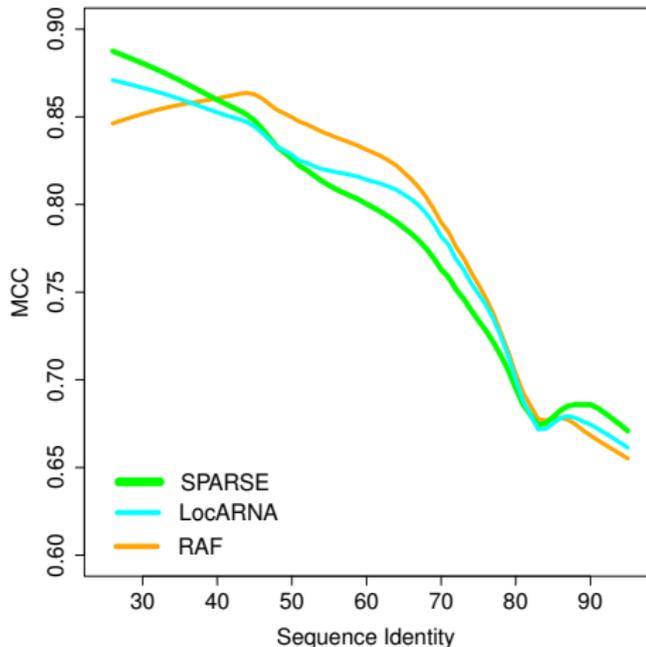
Tool	Sparsification			Mean Time per Instance	Speedup vs. LocARNA
	θ_1	θ_2	θ_3		
LocARNA	1e-3	-	-	2.02s	1.0
SPARSE	1e-3	1e-5	1e-4	0.92s	2.2
RAF	2e-3	-	-	0.37s	5.5

Bralibase 2.1, pairwise alignments

Alignment and Prediction Accuracy (Bralibase 2.1, 3-way alignments)



SPS: alignment quality



MCC: prediction quality

Conclusions

SPARSE: very efficient RNA alignment without sequence-based heuristics

- PARSE is **THE** lightweight Sankoff-variant (cf. PMcomp)
 - predicts deleted/inserted base pairs; like original SAF
- SPARSE = **Sparsified** PARSE
 - Novel ensemble-based sparsification (*in-loop* probabilities)
 - No sequence-based heuristics
 - Speeds up SAF: Quadratic Time [$\leftarrow O(n^6)$]

Conclusions

SPARSE: very efficient RNA alignment without sequence-based heuristics

- PARSE is **THE** lightweight Sankoff-variant (cf. PMcomp)
 - predicts deleted/inserted base pairs; like original SAF
- SPARSE = **Sparsified** PARSE
 - Novel ensemble-based sparsification (*in-loop* probabilities)
 - No sequence-based heuristics
 - Speeds up SAF: Quadratic Time [$\leftarrow O(n^6)$]

Conclusions

SPARSE: very efficient RNA alignment without sequence-based heuristics

- PARSE is **THE** lightweight Sankoff-variant (cf. PMcomp)
 - predicts deleted/inserted base pairs; like original SAF
- SPARSE = **Sparsified** PARSE
 - Novel **ensemble-based sparsification** (*in-loop* probabilities)
 - No sequence-based heuristics
 - Speeds up SAF: **Quadratic Time** [$\leftarrow O(n^6)$]

Conclusions

SPARSE: very efficient RNA alignment without sequence-based heuristics

- PARSE is **THE** lightweight Sankoff-variant (cf. PMcomp)
 - predicts deleted/inserted base pairs; like original SAF
- SPARSE = **Sparsified** PARSE
 - Novel **ensemble-based sparsification** (*in-loop* probabilities)
 - No sequence-based heuristics
 - Speeds up SAF: **Quadratic Time** [$\leftarrow O(n^6)$]

Conclusions

SPARSE: very efficient RNA alignment without sequence-based heuristics

- PARSE is **THE** lightweight Sankoff-variant (cf. PMcomp)
 - predicts deleted/inserted base pairs; like original SAF
- SPARSE = **Sparsified** PARSE
 - Novel **ensemble-based sparsification** (*in-loop* probabilities)
 - No sequence-based heuristics
 - Speeds up SAF: **Quadratic Time** [$\leftarrow O(n^6)$]

Conclusions

SPARSE: very efficient RNA alignment without sequence-based heuristics

- PARSE is **THE** lightweight Sankoff-variant (cf. PMcomp)
 - predicts deleted/inserted base pairs; like original SAF
- SPARSE = **Sparsified** PARSE
 - Novel **ensemble-based sparsification** (*in-loop* probabilities)
 - No sequence-based heuristics
 - Speeds up SAF: **Quadratic Time** [$\leftarrow O(n^6)$]

Thanks

... for your attention

... to my coauthors

- Christina Schmiedl
- Milad Miladi
- Mathias Möhl
- Rolf Backofen



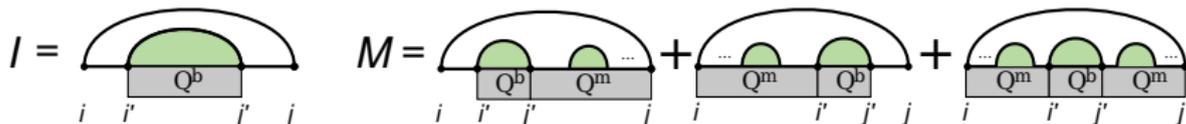
... and the German Research Foundation **DFG**

Appendix

Computing “In Loop” Probabilities

from McCaskill matrices: Q_b, Q_m

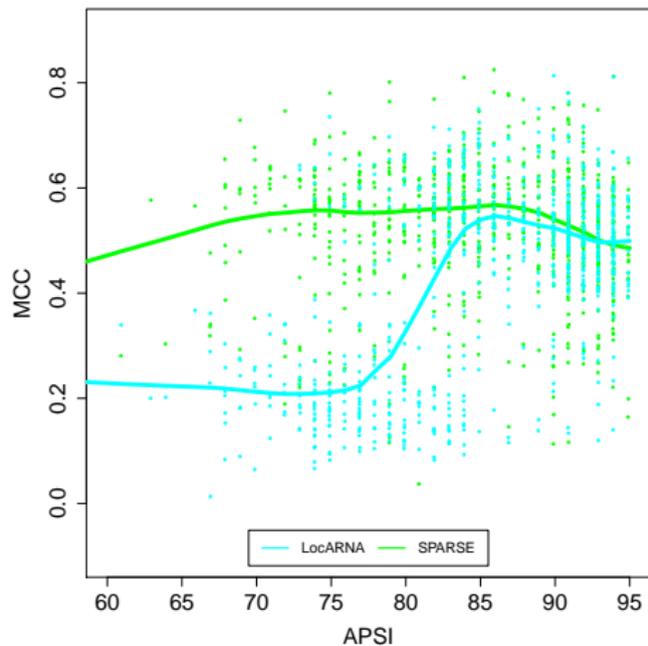
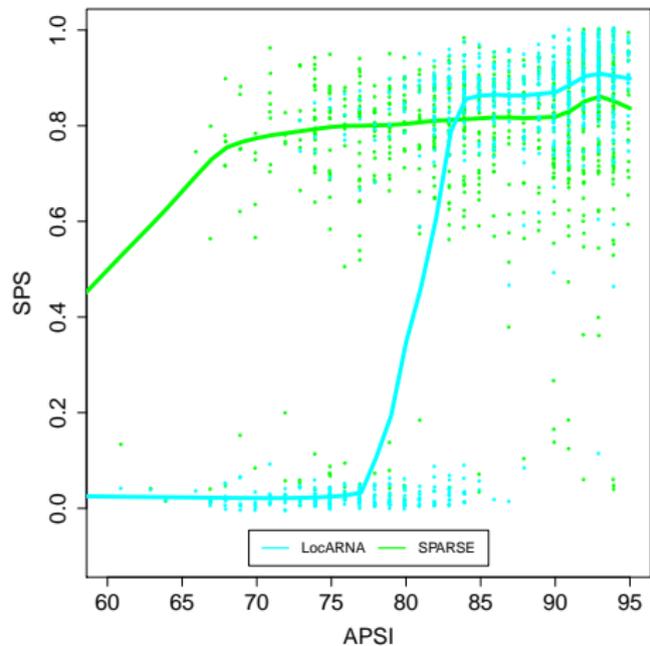
$$\Pr[(i',j') \text{ base pair in loop of } (i,j)] \\ = (I + M)/Q$$



similar: **$\Pr[k \text{ unpaired in loop of } (i,j)]$**

[ExpARNA-P; Schmiedl et al., RECOMB 2012]

SPARSE Improves Over LocARNA for Specific Families



(shown: IRES HCV, pairwise)