

A partition function algorithm for RNA-RNA interaction

Hamidreza Chitsaz

Raheleh Salari, Cenk Sahinalp, Rolf Backofen

Wayne State University
chitsaz@wayne.edu

Benasque RNA Meeting

July 27th, 2012



Mini biography

Robotics → RNA → Genome Assembly

- ▶ **University of Illinois, Urbana-Champaign (Steven M. LaValle): PhD, Computer Science, 2008**
- ▶ Simon Fraser University, Vancouver (Cenk Sahinalp): Postdoc, RNA algorithms, 2009
- ▶ University of California, San Diego (Pavel Pevzner): Postdoc, genome assembly, 2011
- ▶ Wayne State University, Detroit: Assistant professor, 2011-



Mini biography

Robotics → RNA → Genome Assembly

- ▶ University of Illinois, Urbana-Champaign (Steven M. LaValle): PhD, Computer Science, 2008
- ▶ Simon Fraser University, Vancouver (Cenk Sahinalp): Postdoc, RNA algorithms, 2009
- ▶ University of California, San Diego (Pavel Pevzner): Postdoc, genome assembly, 2011
- ▶ Wayne State University, Detroit: Assistant professor, 2011-



Mini biography

Robotics → RNA → Genome Assembly

- ▶ University of Illinois, Urbana-Champaign (Steven M. LaValle): PhD, Computer Science, 2008
- ▶ Simon Fraser University, Vancouver (Cenk Sahinalp): Postdoc, RNA algorithms, 2009
- ▶ University of California, San Diego (Pavel Pevzner): Postdoc, genome assembly, 2011
- ▶ Wayne State University, Detroit: Assistant professor, 2011-



Mini biography

Robotics → RNA → Genome Assembly

- ▶ University of Illinois, Urbana-Champaign (Steven M. LaValle): PhD, Computer Science, 2008
- ▶ Simon Fraser University, Vancouver (Cenk Sahinalp): Postdoc, RNA algorithms, 2009
- ▶ University of California, San Diego (Pavel Pevzner): Postdoc, genome assembly, 2011
- ▶ Wayne State University, Detroit: Assistant professor, 2011-



Efficient *de novo* assembly of single-cell bacterial genomes from short-read data sets

Hamidreza Chitsaz^{1,6}, Joyclyn L Yee-Greenbaum^{2,6}, Glenn Tesler³, Mary-Jane Lombardo², Christopher L Dupont², Jonathan H Badger², Mark Novotny², Douglas B Rusch⁴, Louise J Fraser⁵, Niall A Gormley⁵, Ole Schulz-Trieglaff⁵, Geoffrey P Smith⁵, Dirk J Evers⁵, Pavel A Pevzner¹ & Roger S Lasken²

Whole genome amplification by the multiple displacement amplification (MDA) method allows sequencing of DNA from single cells of bacteria that cannot be cultured. Assembling a genome is challenging, however, because MDA generates highly

BIOINFORMATICS

Vol. 28 ISMB 2012, pages 1188–1196
doi:10.1093/bioinformatics/bts219

SEQuel: improving the accuracy of genome assemblies

Roy Ronen^{1,†}, Christina Boucher^{2,†}, Hamidreza Chitsaz³ and Pavel Pevzner^{2,*}

¹Bioinformatics Graduate Program, ²Department of Computer Science and Engineering, University of California, San Diego, La Jolla, CA 92093 and ³Department of Computer Science, Wayne State University, Detroit, MI 48202, USA

ABSTRACT

Motivation: Assemblies of next-generation sequencing (NGS) data,

finished genomes assembled using the previous technologies (Alkan, *et al.*, 2011). Earlier assembly algorithms developed for Sanger



Algorithmic Biology Laboratory

Wayne State University

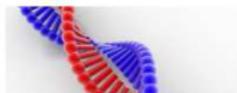


Algorithmic Biology Laboratory

[Home](#)
[Publications](#)
[Research](#)
[Software](#)
[People](#)
[Links](#)
[Contact Us](#)



Welcome to Algorithmic Biology Laboratory:
ABL research focuses on quantitative modeling of various life phenomena at cellular and molecular levels. In particular, ABL research includes modeling RNA structure and ncRNA-mRNA interactions, *de novo* assembly of short and long sequencing



NEWS

Single Cell Genomics
Our single cell assembly paper entitled "*Efficient de novo assembly of single-cell bacterial genomes from short-read data sets*" has been published in [Nature Biotechnology](#). We modified Velvet to obtain Velvet-SC for MDA single cell data.



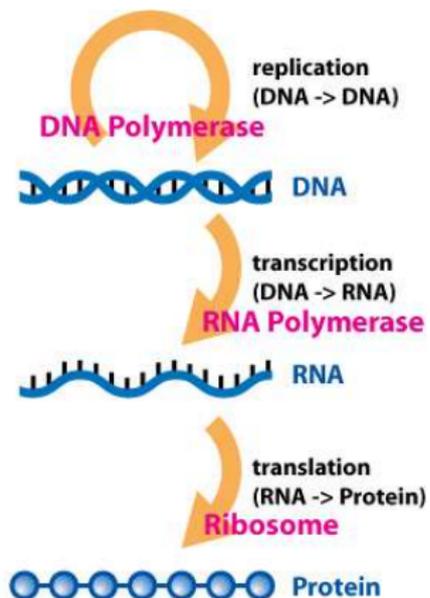
Camel Released
We have released Camel, a tool for correcting substitution errors in Illumina reads, particularly

<http://compbio.cs.wayne.edu>



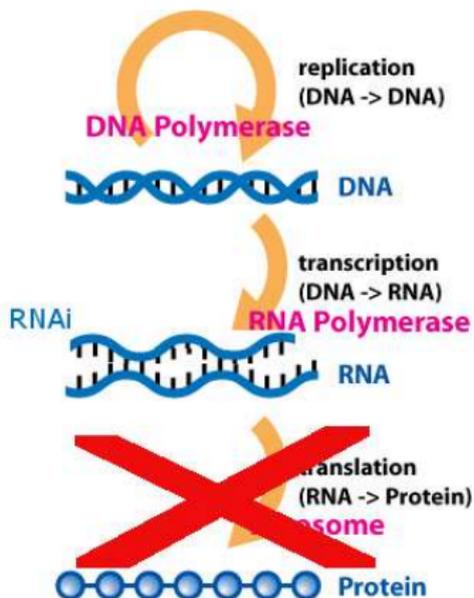
Central dogma

DNA → RNA → Protein



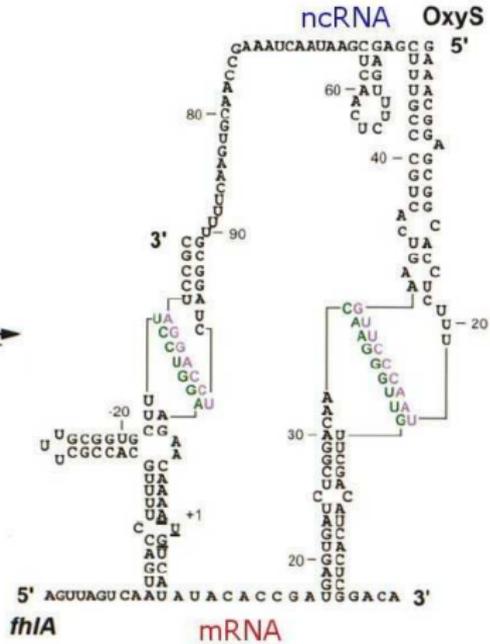
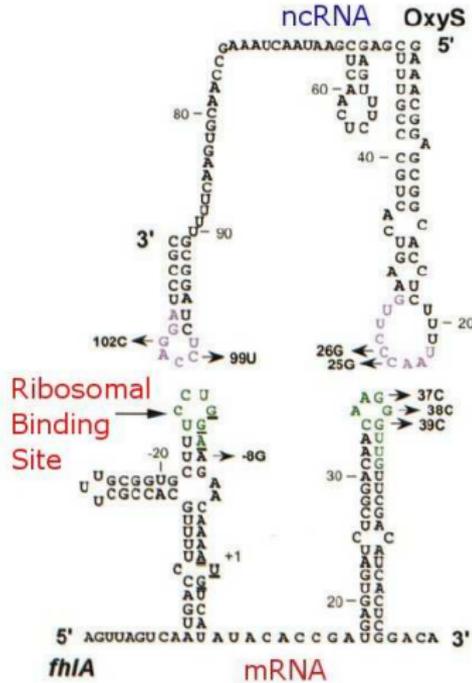
Motivation

Post-transcriptional regulation of gene expression



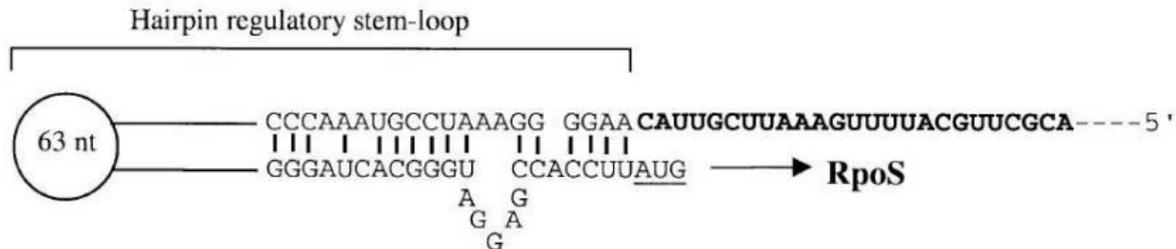
Regulatory RNA

Repression example (Argaman and Altuvia, J. Mol. Biol. 2000)

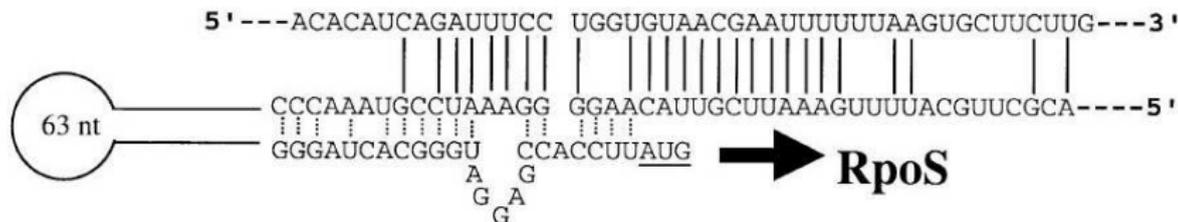


Regulatory RNA

Activation example (Repoila, Majdalani, and Gottesman, Mol. Microbiol. 2003)



DsrA



Background

RNA-RNA MFE structure prediction

- ▶ **Avoid intramolecular base pairing**
RNAhybrid (Rehmsmeier *et al.* 2004), **RNA duplex** (Bernhart *et al.* 2006), **UNAFold** (Markham *et al.* 2008)
No internal structure
- ▶ Concatenate input sequences as a single strand; no pseudoknots
PairFold (Andronescu *et al.* 2005), **RNAcofold** (Bernhart *et al.* 2006)
No kissing hairpins
- ▶ Predict binding sites
RNAup (Mückstein *et al.* 2008), **intaRNA** (Busch *et al.* 2008)
Just one binding site not complete structure
- ▶ Concatenate input sequences; consider special pseudoknots
NUPACK (Dirks *et al.* 2003,2007)
Still no kissing hairpins!



Background

RNA-RNA MFE structure prediction

- ▶ Avoid intramolecular base pairing
RNAhybrid (Rehmsmeier *et al.* 2004), **RNA duplex** (Bernhart *et al.* 2006), **UNAFold** (Markham *et al.* 2008)
No internal structure
- ▶ Concatenate input sequences as a single strand; no pseudoknots
PairFold (Andronescu *et al.* 2005), **RNAcofold** (Bernhart *et al.* 2006)
No kissing hairpins
- ▶ Predict binding sites
RNAup (Mückstein *et al.* 2008), **intaRNA** (Busch *et al.* 2008)
Just one binding site not complete structure
- ▶ Concatenate input sequences; consider special pseudoknots
NUPACK (Dirks *et al.* 2003,2007)
Still no kissing hairpins!



Background

RNA-RNA MFE structure prediction

- ▶ Avoid intramolecular base pairing
RNAhybrid (Rehmsmeier *et al.* 2004), **RNA duplex** (Bernhart *et al.* 2006), **UNAFold** (Markham *et al.* 2008)
No internal structure
- ▶ Concatenate input sequences as a single strand; no pseudoknots
PairFold (Andronescu *et al.* 2005), **RNAcofold** (Bernhart *et al.* 2006)
No kissing hairpins
- ▶ Predict binding sites
RNAup (Mückstein *et al.* 2008), **intaRNA** (Busch *et al.* 2008)
Just one binding site not complete structure
- ▶ Concatenate input sequences; consider special pseudoknots
NUPACK (Dirks *et al.* 2003,2007)
Still no kissing hairpins!



Background

RNA-RNA MFE structure prediction

- ▶ Avoid intramolecular base pairing
RNAhybrid (Rehmsmeier *et al.* 2004), **RNA duplex** (Bernhart *et al.* 2006), **UNAFold** (Markham *et al.* 2008)
No internal structure
- ▶ Concatenate input sequences as a single strand; no pseudoknots
PairFold (Andronescu *et al.* 2005), **RNAcofold** (Bernhart *et al.* 2006)
No kissing hairpins
- ▶ Predict binding sites
RNAup (Mückstein *et al.* 2008), **intaRNA** (Busch *et al.* 2008)
Just one binding site not complete structure
- ▶ Concatenate input sequences; consider special pseudoknots
NUPACK (Dirks *et al.* 2003,2007)
Still no kissing hairpins!



Background (continued)

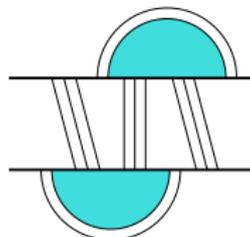
RNA-RNA MFE structure prediction

Consider inter- and intramolecular base pairing

IRIS (Pervouchine 2004), **inteRNA** (Alkan *et al.* 2005), **Grammatical Approach** (Kato *et al.* 2009)

Voilà, now we are talking business.

The problem is NP-Hard (Alkan *et al.* 2005); no surprise as pseudoknots are NP-Hard. Exclude *zigzags* and crossing interactions to lift the curse of complexity and obtain an exact $O(n^6)$ -time $O(n^4)$ -space DP algorithm (albeit for simple base-pair counting).



First order zigzag. A general zigzag involves an arbitrary number of kissing hairpins.



Ahhh...but MFE is often wrong!

Question: how about

1. computing base pairing probabilities,
2. sampling from the Boltzmann ensemble of interaction structures, clustering, centroids, etc.,
3. and computing equilibrium concentrations and melting temperature for RNA-RNA compounds?

Answer: the key enabling technology is the **partition function**. All of the above can be computed from the partition function.



Ahhh...but MFE is often wrong!

Question: how about

1. computing base pairing probabilities,
2. sampling from the Boltzmann ensemble of interaction structures, clustering, centroids, etc.,
3. and computing equilibrium concentrations and melting temperature for RNA-RNA compounds?

Answer: the key enabling technology is the **partition function**. All of the above can be computed from the partition function.



Ahhh...but MFE is often wrong!

Question: how about

1. computing base pairing probabilities,
2. sampling from the Boltzmann ensemble of interaction structures, clustering, centroids, etc.,
3. and computing equilibrium concentrations and melting temperature for RNA-RNA compounds?

Answer: the key enabling technology is the **partition function**. All of the above can be computed from the partition function.



Ahhh...but MFE is often wrong!

Question: how about

1. computing base pairing probabilities,
2. sampling from the Boltzmann ensemble of interaction structures, clustering, centroids, etc.,
3. and computing equilibrium concentrations and melting temperature for RNA-RNA compounds?

Answer: the key enabling technology is the **partition function**. All of the above can be computed from the partition function.



Partition function

$$Q(T) = \sum_{s \in S} e^{-G_s/RT},$$

S = All considered interaction structures,

$$p(s) \propto e^{-G_s/RT},$$

and Q is the normalizing factor. Also other thermodynamic quantities can be derived from Q.



Partition function

$$Q(T) = \sum_{s \in S} e^{-G_s/RT},$$

S = All considered interaction structures,

$$p(s) \propto e^{-G_s/RT},$$

and Q is the normalizing factor. Also other thermodynamic quantities can be derived from Q.



Partition function hardness \geq MFE hardness

Partition function

$$\sum_{s \in S} e^{-G_s/RT}.$$

MFE secondary structure

$$\operatorname{argmin}_{s \in S} G_s.$$

Transform a partition function algorithm to an MFE algorithm by

$$e^{-G_s} \rightarrow G_s$$

$$\times \rightarrow +$$

$$\Sigma \rightarrow \min$$



Partition function hardness \geq MFE hardness

Partition function

$$\sum_{s \in S} e^{-G_s/RT}.$$

MFE secondary structure

$$\operatorname{argmin}_{s \in S} G_s.$$

Transform a partition function algorithm to an MFE algorithm by

$$e^{-G_s} \rightarrow G_s$$

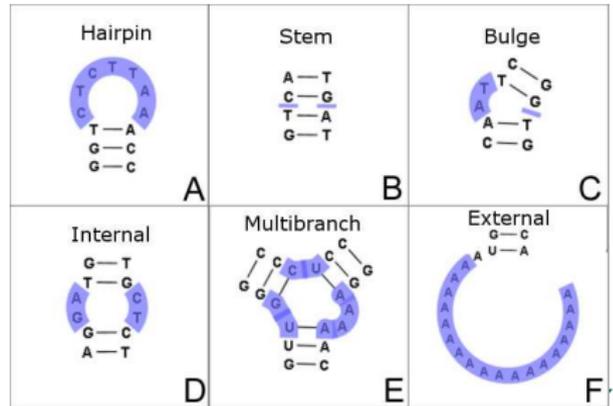
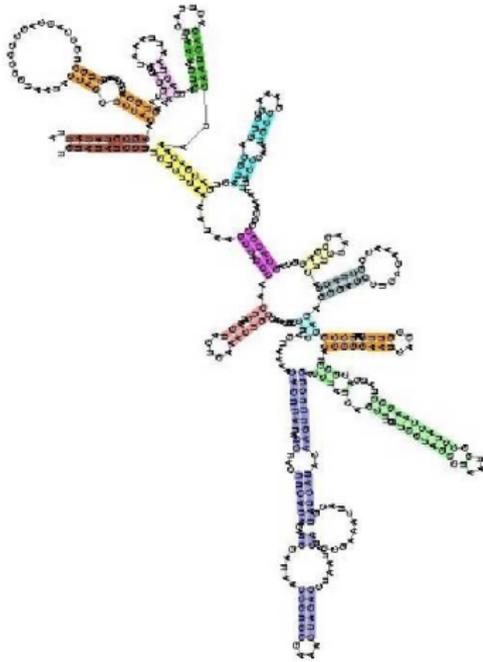
$$\times \rightarrow +$$

$$\Sigma \rightarrow \min$$



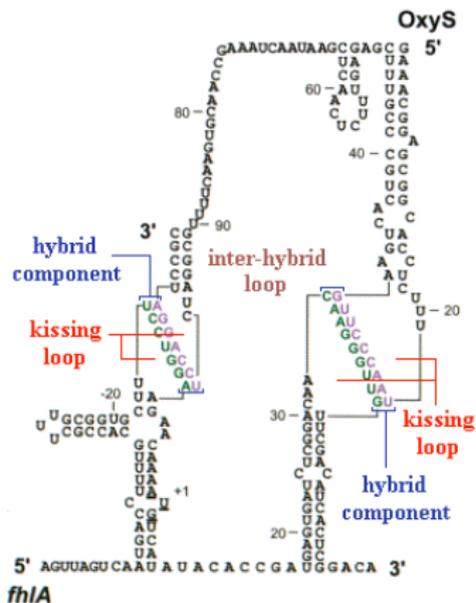
Turner energy model

Mathews *et al.* 1999



Our extension of the Turner model

Chitsaz *et al.*, *Bioinformatics* 25(12): i365-i373



Hybrid component: as if intramolecular, with penalties.

Kissing loop: like multibranch loop.



Interaction partition function

How?

Divide and conquer using dynamic programming:

$$\begin{aligned} Q(T) &= \sum_{s \in S} e^{-G_s/RT} \\ &= \sum_{s=s_a \cup s_b} e^{-(G_{s_a}+G_{s_b})/RT} \\ &= \left[\sum_{s_a \in S_a} e^{-G_{s_a}/RT} \right] \left[\sum_{s_b \in S_b} e^{-G_{s_b}/RT} \right] \\ &= Q_a(T) Q_b(T). \end{aligned}$$

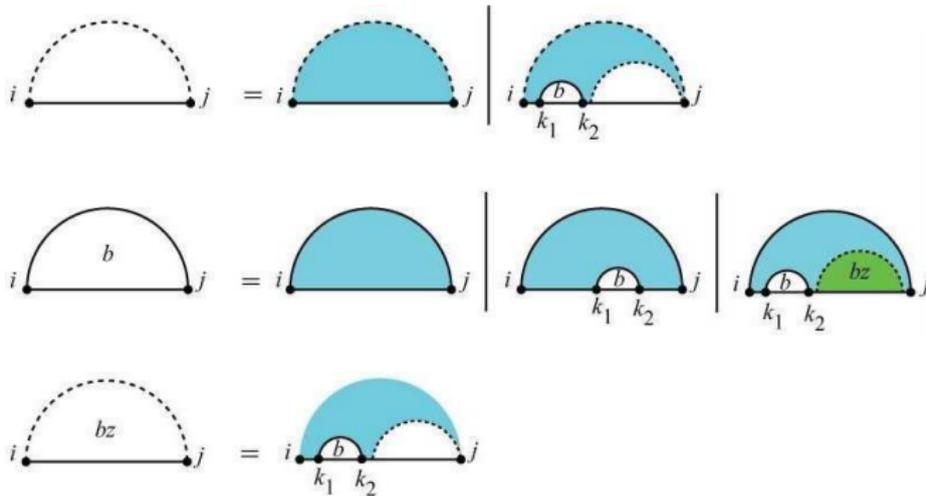


Partition function for single strand (McCaskill 1990)

straight horizontal line: nucleotides indexed from 1 to n

solid arc: a base pair

dashed arc: can be base pair or not



white region: open to more recursions

blue region: finalized in the recursion, compute its energy contribution

green region: open to more recursions with multibranch loop energy

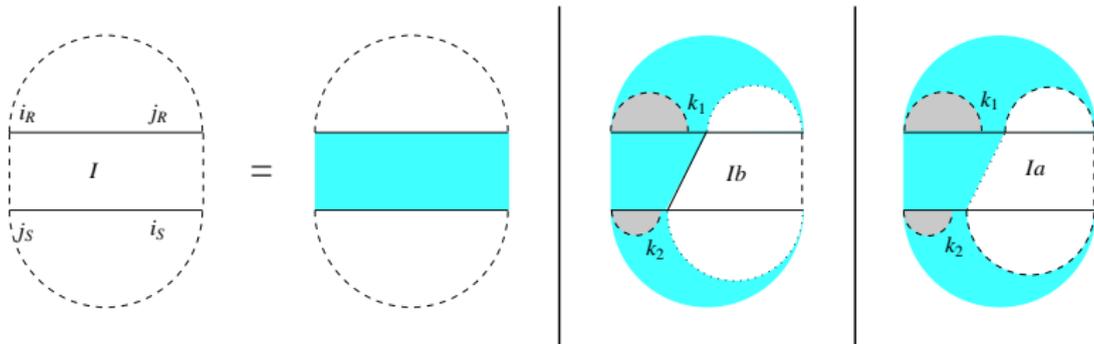
Partition function for two strands

straight vertical line: intermolecular bond

solid: a base pair

dotted: not a base pair

dashed: either of those two



$$Q_{i_R, j_R, i_S, j_S}^I = Q_{i_R, j_R} Q_{i_S, j_S} + \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1-1} Q_{k_2+1, j_S} Q_{k_1, j_R, i_S, k_2}^{Ib} + \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1-1} Q_{k_2+1, j_S} Q_{k_1, j_R, i_S, k_2}^{Ia}$$

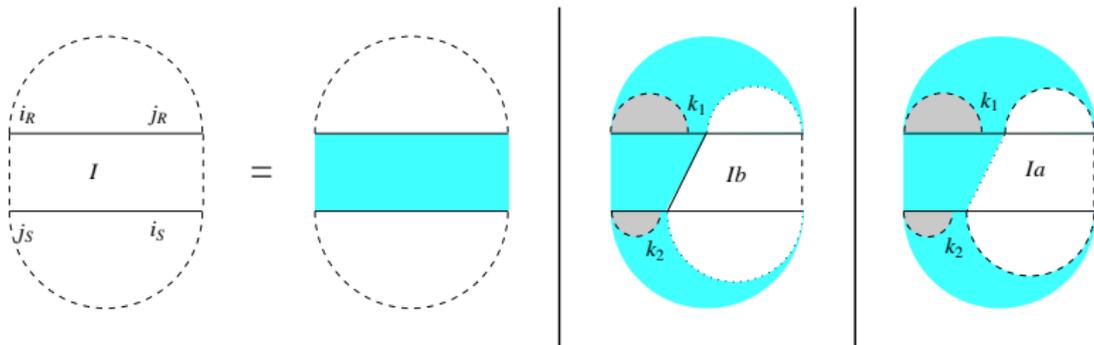
Partition function for two strands

straight vertical line: intermolecular bond

solid: a base pair

dotted: not a base pair

dashed: either of those two



$$Q_{i_R, j_R, i_S, j_S}^I = Q_{i_R, j_R} Q_{i_S, j_S} + \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1 - 1} Q_{k_2 + 1, j_S} Q_{k_1, j_R, i_S, k_2}^{Ib}$$

$$+ \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1 - 1} Q_{k_2 + 1, j_S} Q_{k_1, j_R, i_S, k_2}^{Ia}$$

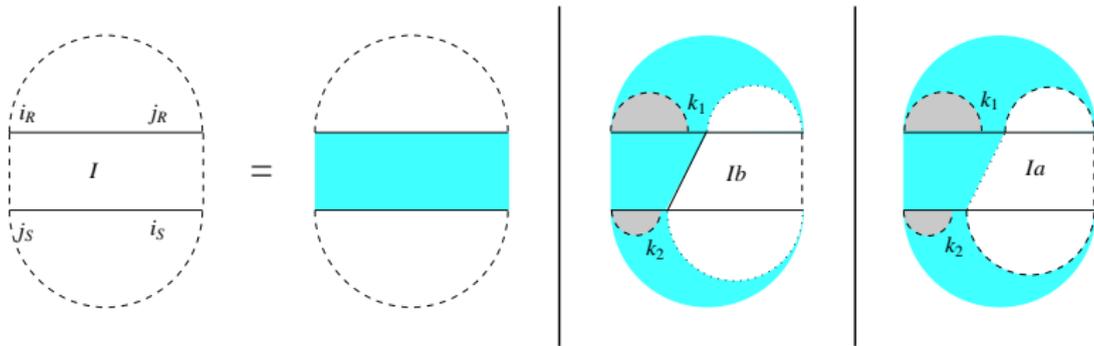
Partition function for two strands

straight vertical line: intermolecular bond

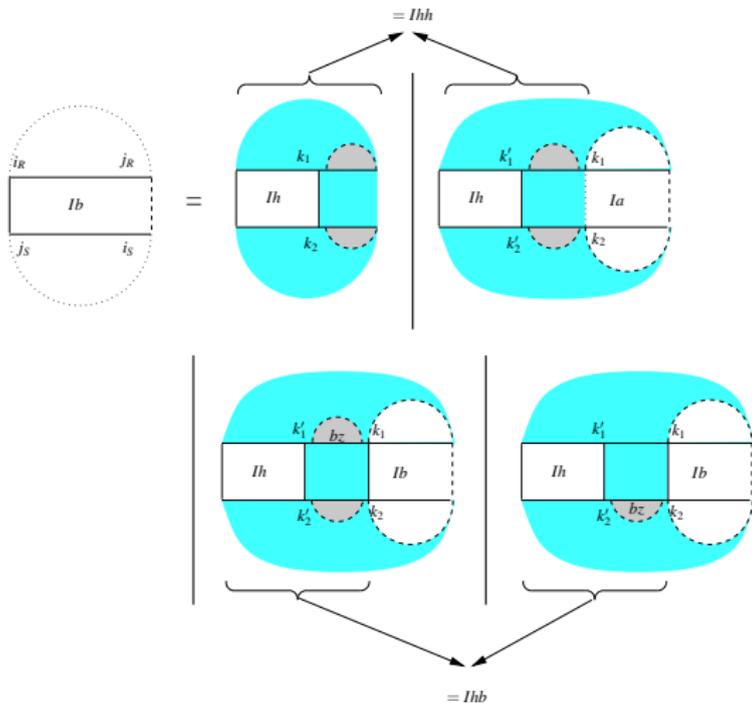
solid: a base pair

dotted: not a base pair

dashed: either of those two

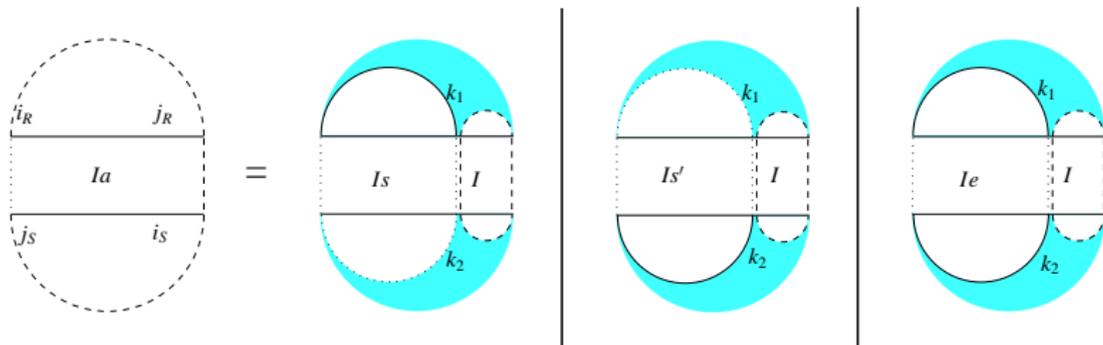


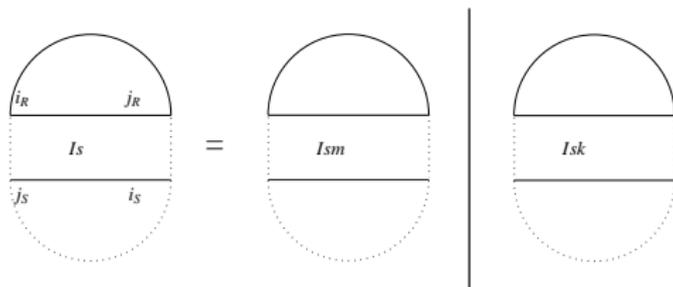
$$Q_{i_R, j_R, i_S, j_S}^I = Q_{i_R, j_R} Q_{i_S, j_S} + \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1-1} Q_{k_2+1, j_S} Q_{k_1, j_R, i_S, k_2}^{Ib} + \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1-1} Q_{k_2+1, j_S} Q_{k_1, j_R, i_S, k_2}^{Ia}$$



b: stands for bond

a: stands for arc
 s: stands for subsume
 e: stands for equivalent

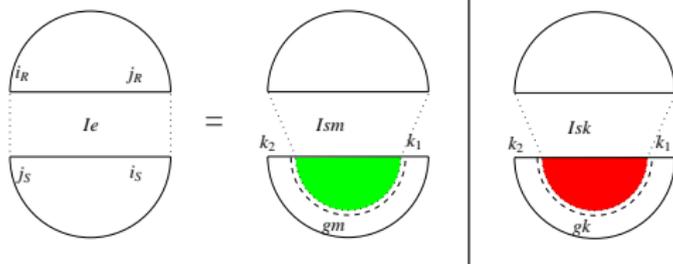




s : stands for subsume

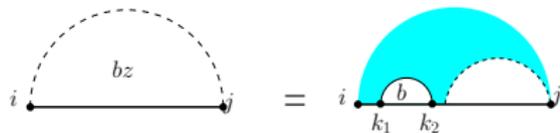
k : stands for kissing-loop

m : stands for multi-loop

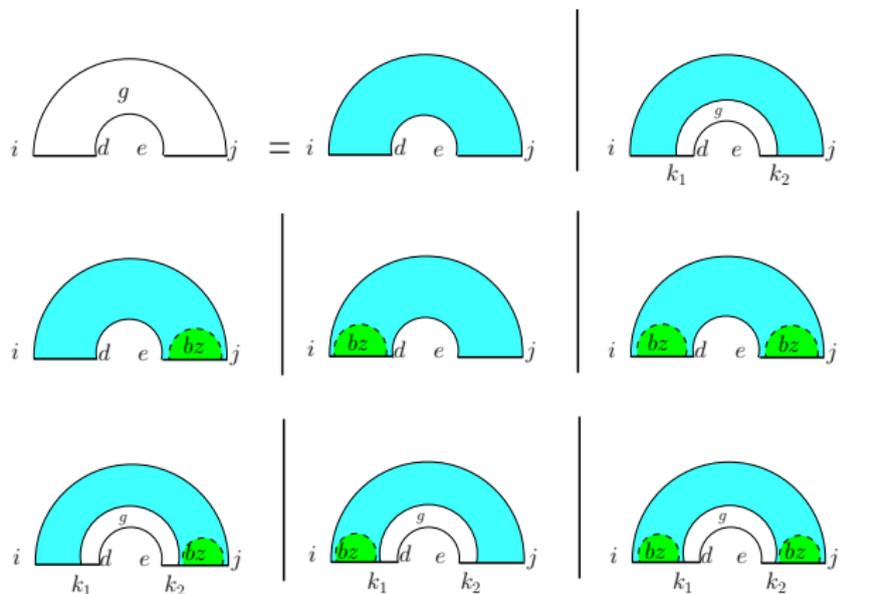


e : stands for equivalent

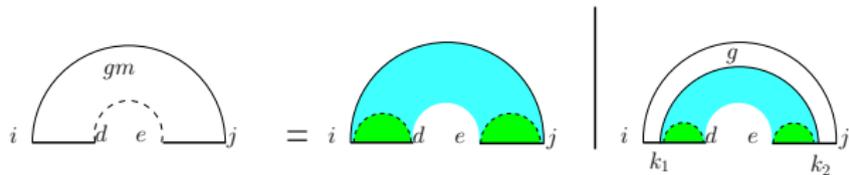
All tables



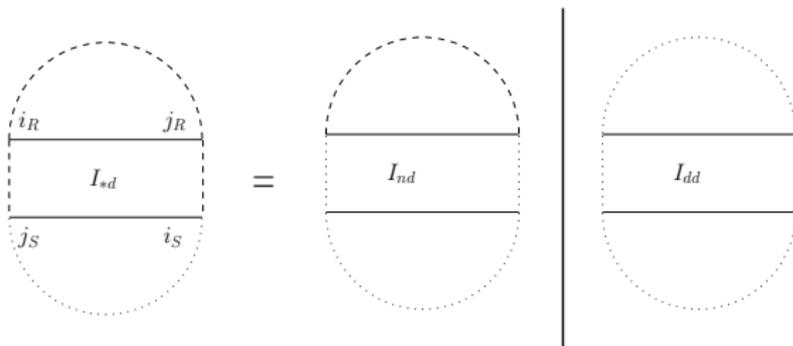
All tables



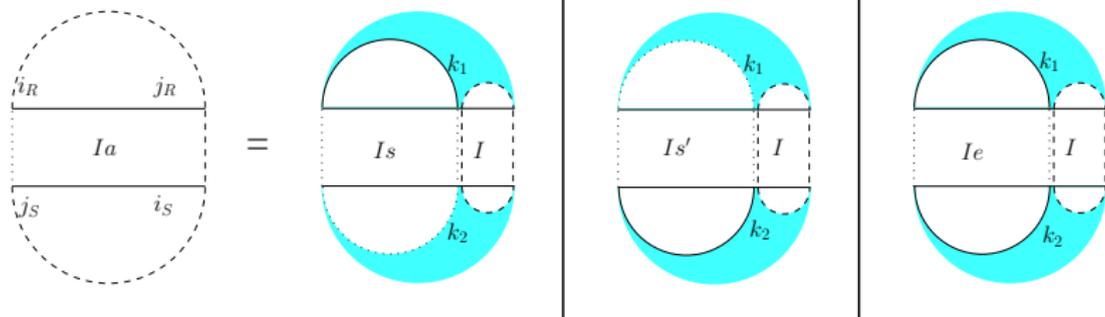
All tables



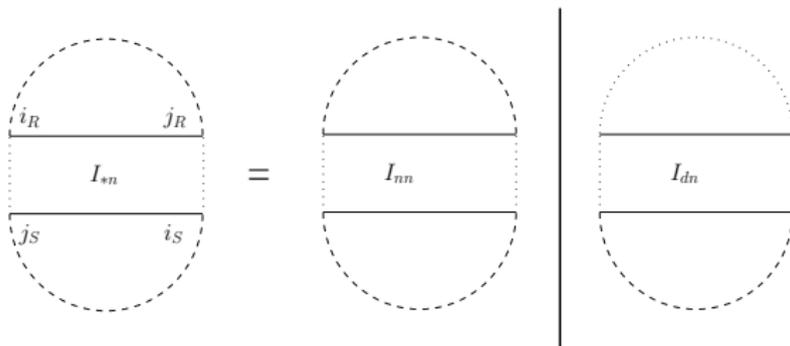
All tables



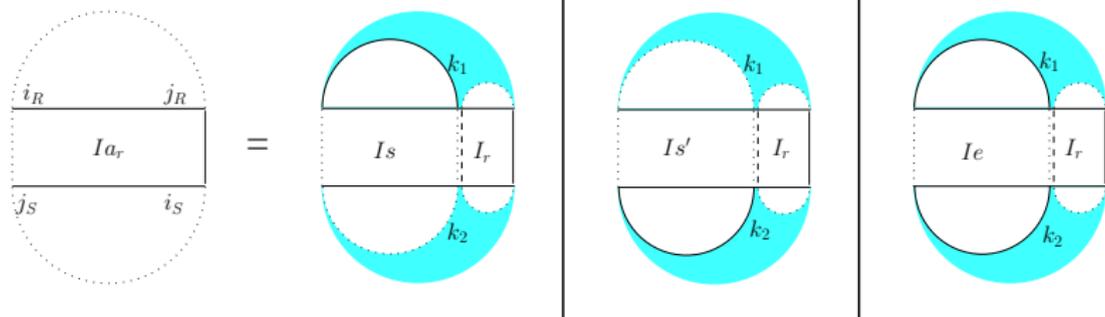
All tables



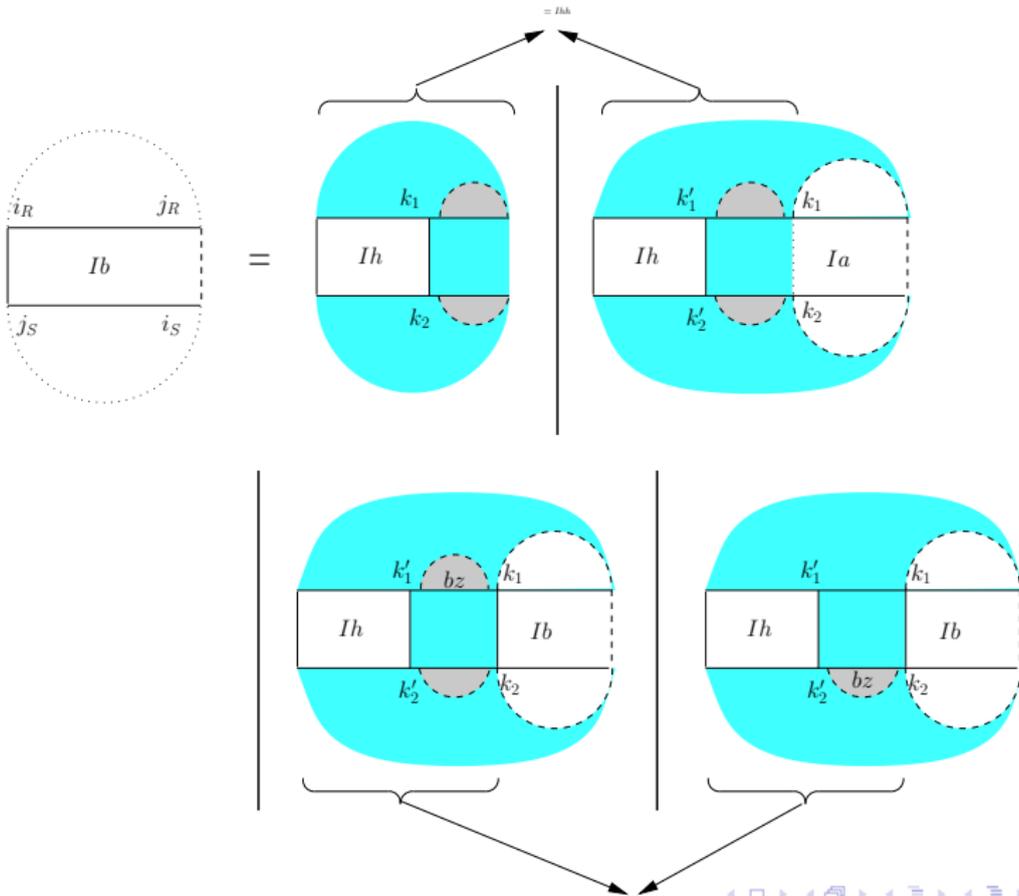
All tables



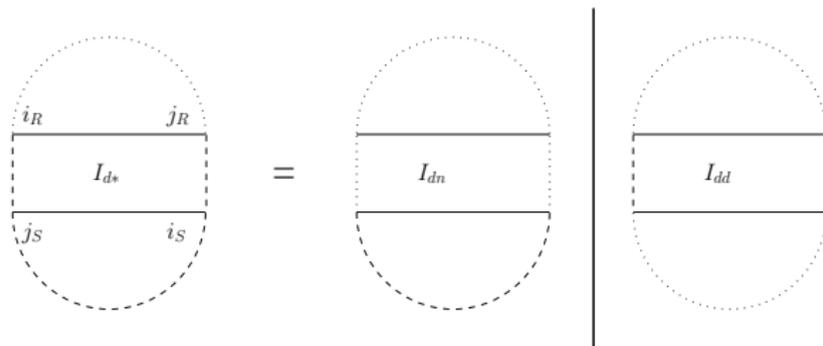
All tables



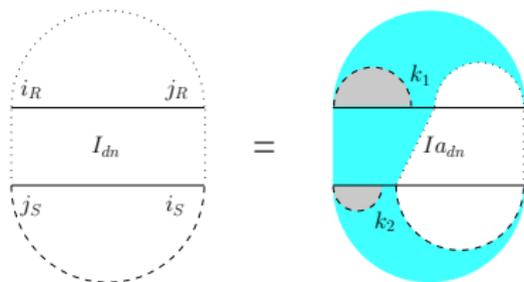
All tables



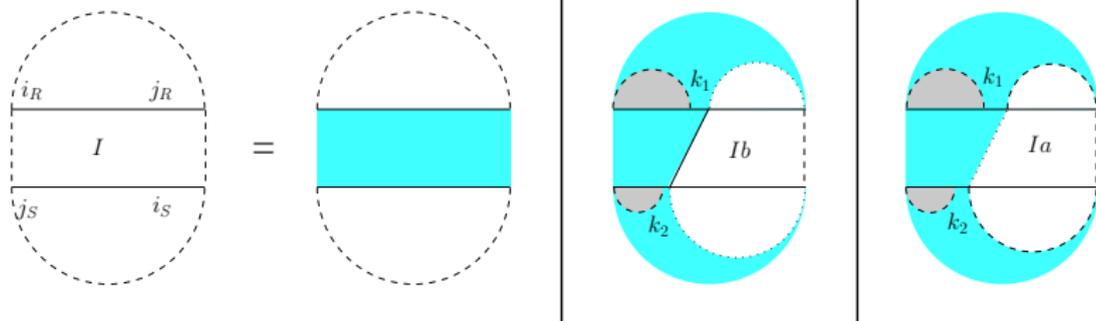
All tables



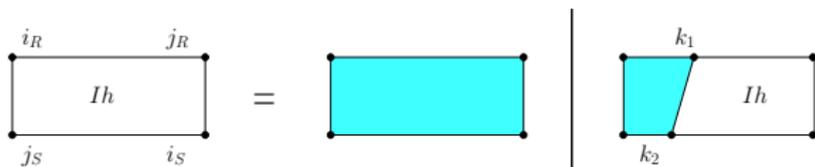
All tables



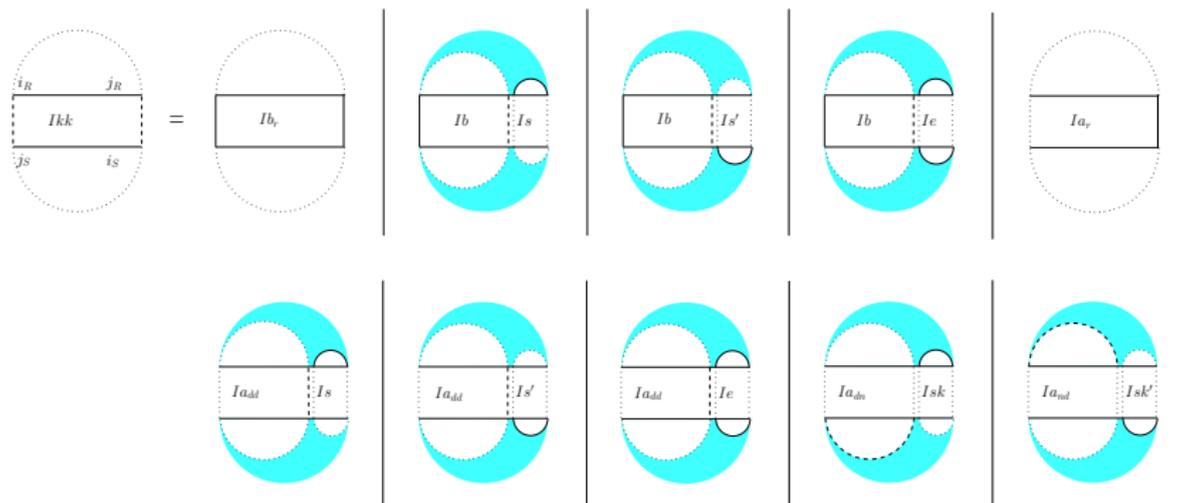
All tables



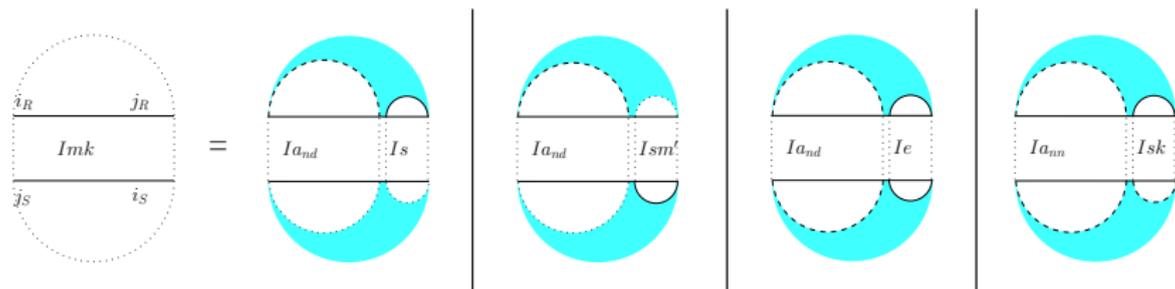
All tables



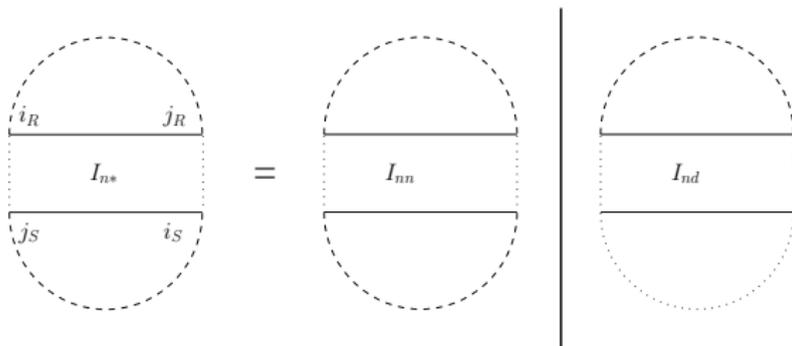
All tables



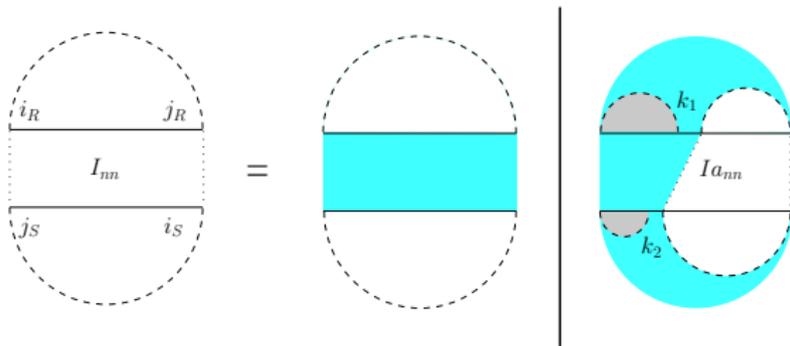
All tables



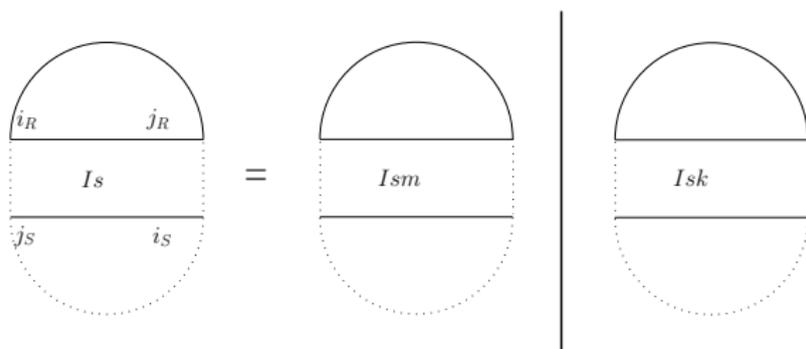
All tables



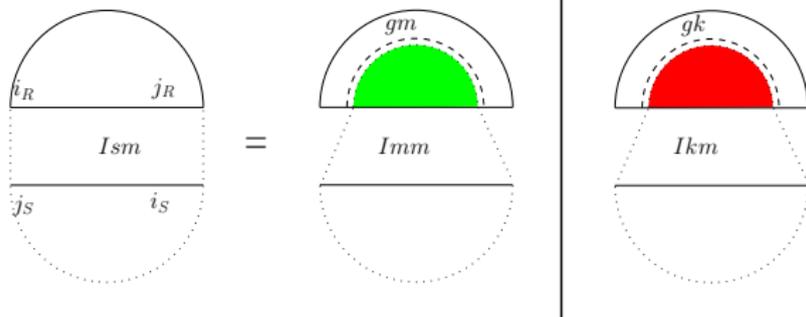
All tables



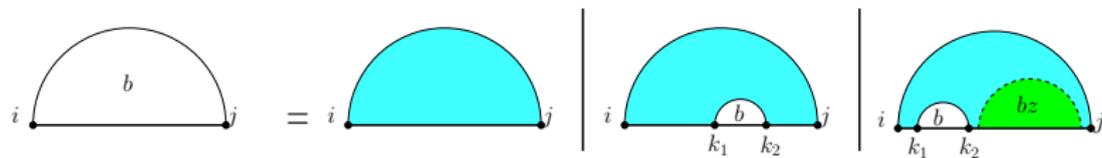
All tables



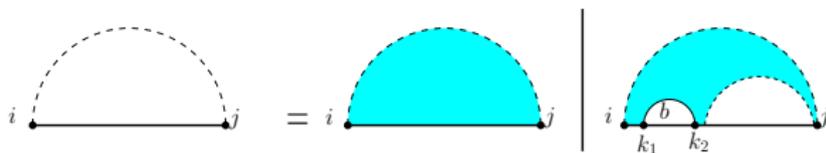
All tables



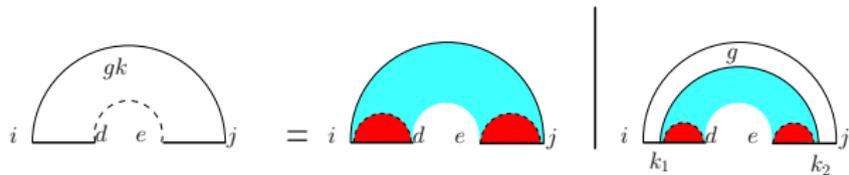
All tables



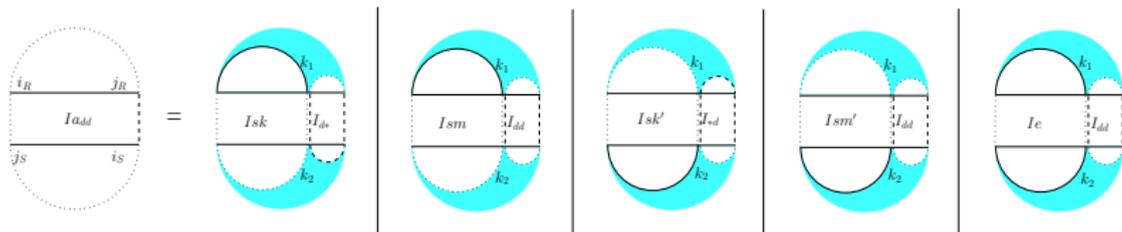
All tables



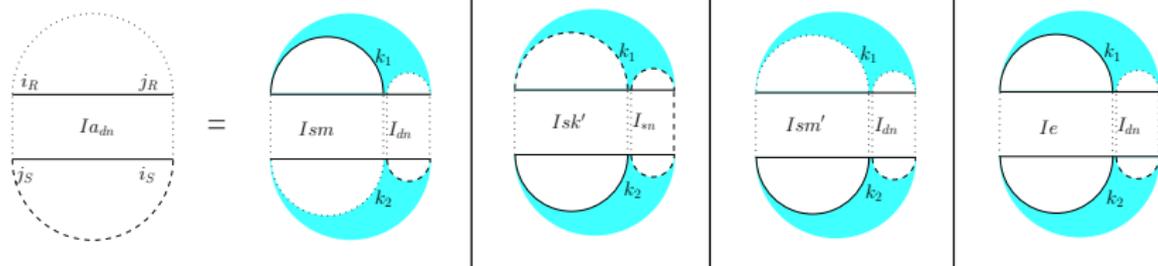
All tables



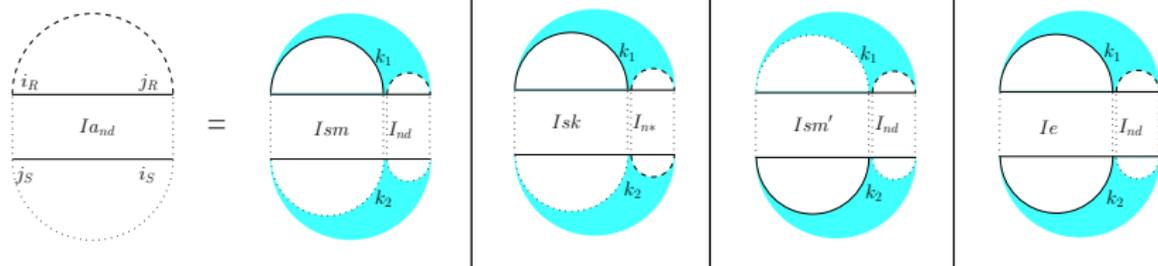
All tables



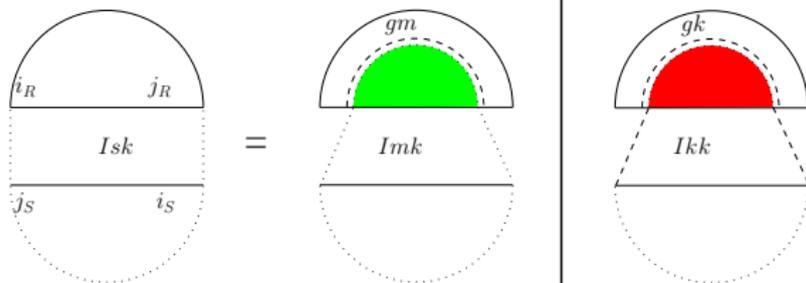
All tables



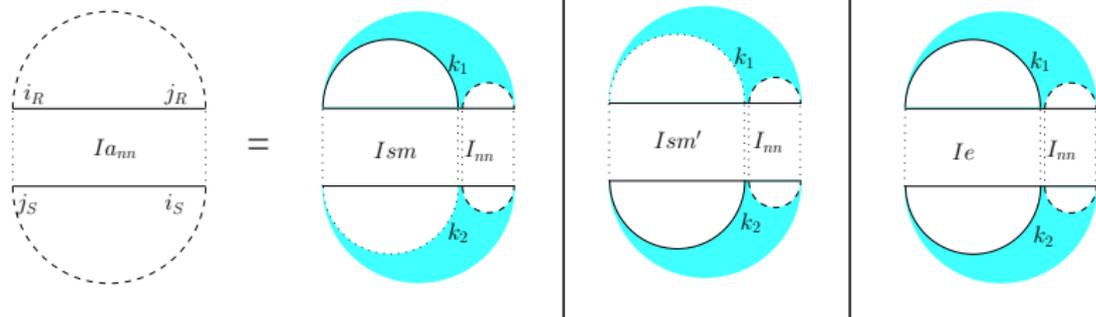
All tables



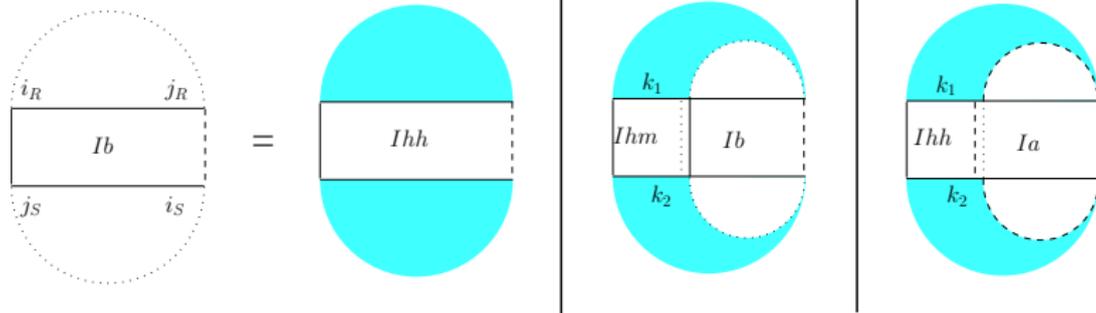
All tables



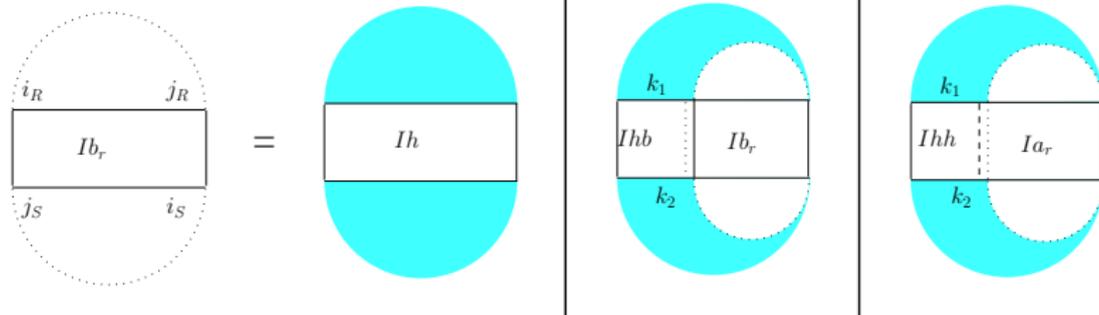
All tables



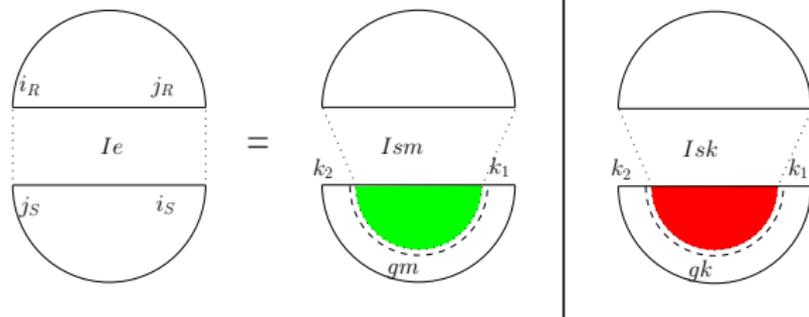
All tables



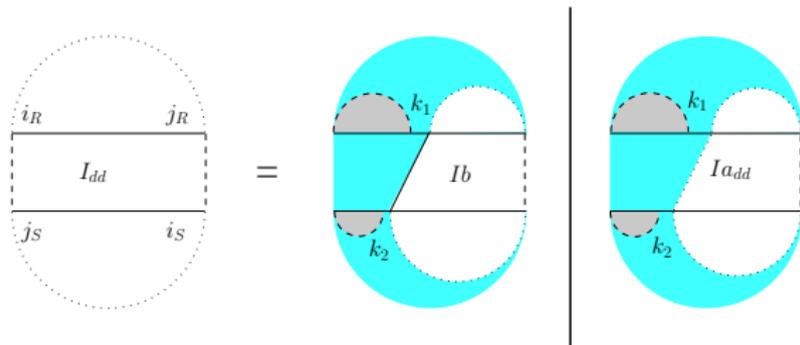
All tables



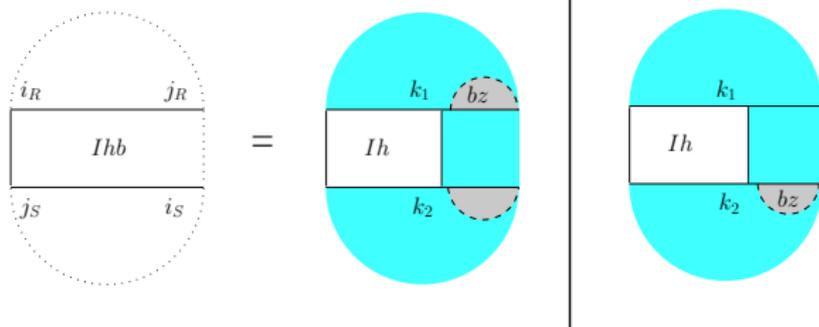
All tables



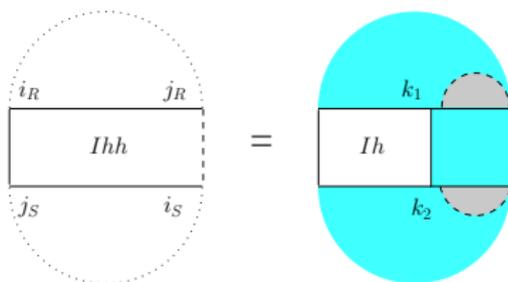
All tables



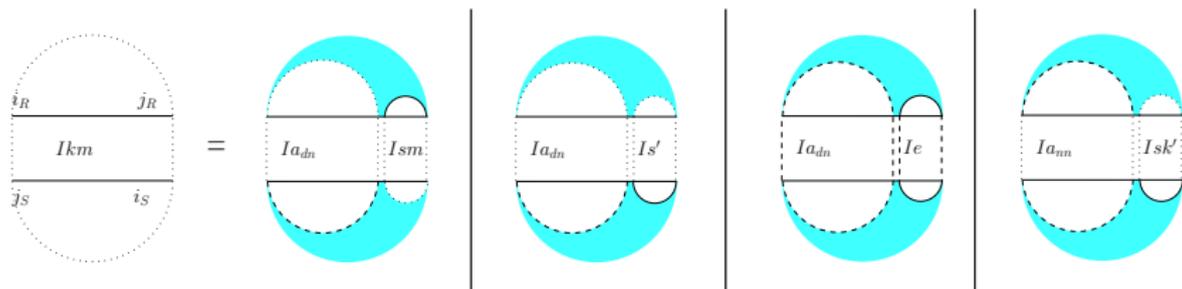
All tables



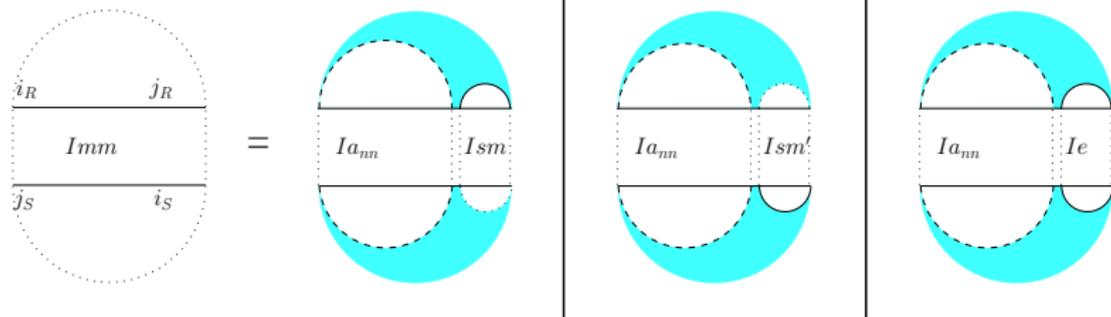
All tables



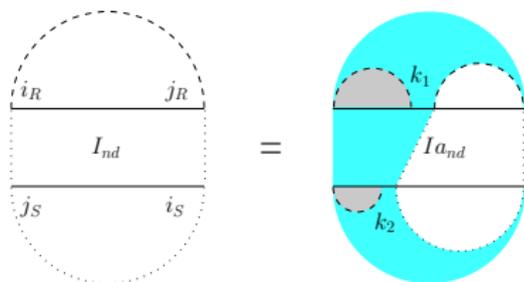
All tables



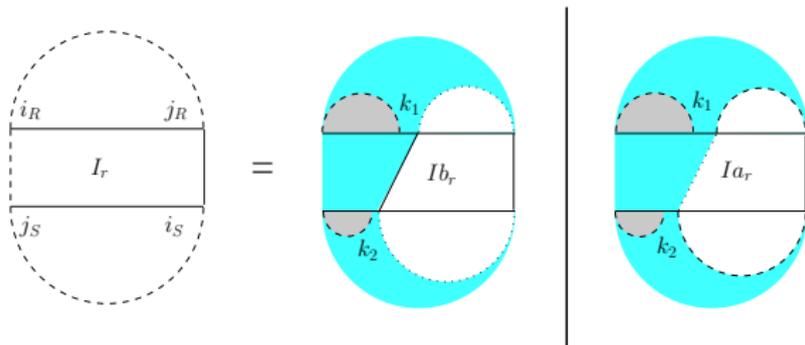
All tables



All tables



All tables



Equilibrium concentrations

For two RNAs **R** and **S**

Assume five types of chemical compounds: **R**, **S**, **RR**, **SS**, **RS**.

Solve

$$K_R = \frac{Q'_{RR}}{Q_R^2} = \frac{N_{RR}}{N_R^2},$$

$$K_S = \frac{Q'_{SS}}{Q_S^2} = \frac{N_{SS}}{N_S^2},$$

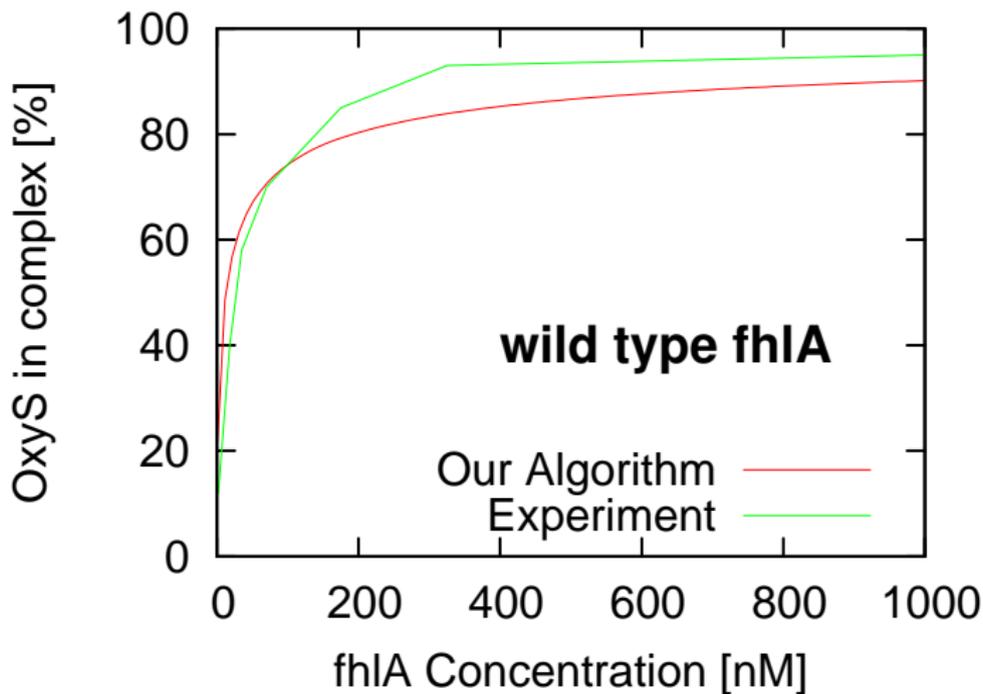
$$K_{RS} = \frac{Q'_{RS}}{Q_R Q_S} = \frac{N_{RS}}{N_R N_S},$$

$$N_{RS} = N_R^0 - 2N_{RR} - N_S = N_S^0 - 2N_{SS} - N_R,$$

to obtain the equilibrium concentrations N . N^0 are the initial concentrations of single strands.



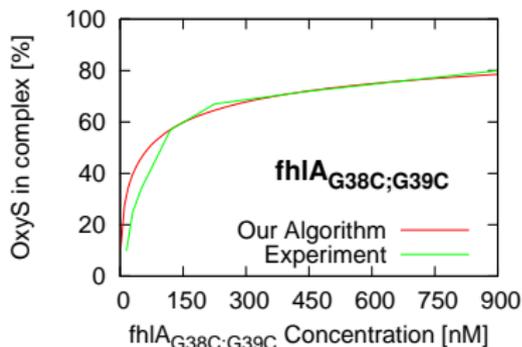
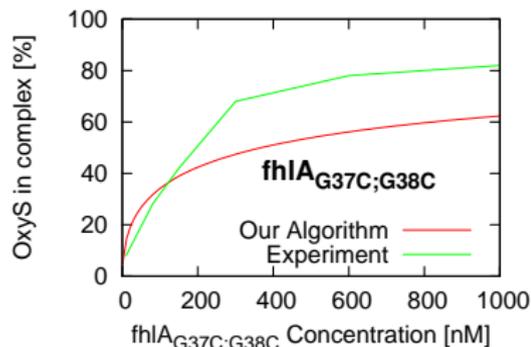
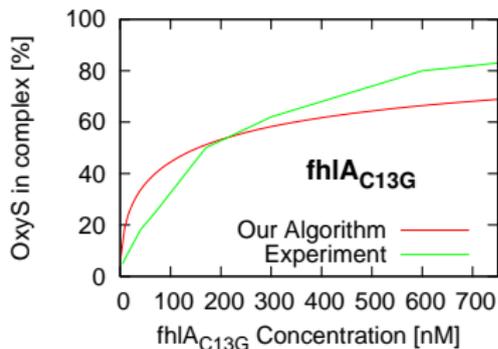
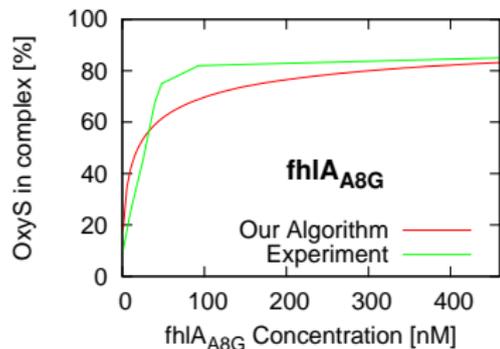
Equilibrium concentration of OxyS with wild type fhIA



Init. [OxyS] = 2nM, [fhIA] = 0 to 1000nM



Equilibrium concentration of OxyS with fhIA mutants



Melting temperature prediction

Comparison of piRNA results over three data sets

Set	Size	Length	Avg error		
			piRNA	RNAcofold	UNAFold
I	9 short pairs	5-7nt	1.48°C	9.35°C	8.55°C
II	12 pairs	~ 20nt	4.86°C	22.97°C	9.12°C
III	62 pairs	22 – 40nt	1.91°C	14.34°C	26.53°C

Set	Size	Length	Spearman rank correlation		
			piRNA	RNAcofold	UNAFold
I	9 short pairs	5-7nt	0.97	0.97	0.57
II	12 pairs	~ 20nt	0.41	-0.03	0.1
III	62 pairs	22 – 40nt	0.3	-0.04	0.24



Promised base pairing probabilities

P^l and P^{la} examples

$$P_{i_R, j_R, i_S, j_S}^l = \sum_{\substack{1 \leq k_1 < i_R \\ i_S < k_2 \leq L_S}} P_{k_1, j_R, i_S, k_2}^{la} \frac{(Q_{k_1, i_R, j_S, k_2}^{ls} + Q_{k_1, i_R, j_S, k_2}^{ls'} + Q_{k_1, i_R, j_S, k_2}^{le}) Q_{i_R, j_R, i_S, j_S}^l}{Q_{k_1, j_R, i_S, k_2}^{la}},$$

$$P_{i_R, j_R, i_S, j_S}^{la} = \sum_{\substack{1 \leq k_1 \leq i_R \\ i_S \leq k_2 \leq L_S}} P_{k_1, j_R, i_S, k_2}^l \frac{Q_{k_1, i_R-1} Q_{j_S+1, k_2} Q_{i_R, j_R, i_S, j_S}^{la}}{Q_{k_1, j_R, i_S, k_2}^l} +$$
$$\sum_{\substack{1 \leq k_1 < i_R \\ i_S \leq k_2 \leq L_S}} P_{k_1, j_R, i_S, k_2}^{lb} \frac{Q_{k_1, i_R, j_S, k_2}^{lhh} Q_{i_R, j_R, i_S, j_S}^{la}}{Q_{k_1, j_R, i_S, k_2}^{lb}}.$$

More on this part will be presented by Peter Stadler.



Sampling from the Boltzmann ensemble

- ▶ Push $l(1, n, 1, m)$ onto the stack.
- ▶ Iterate until the stack is empty, i.e. reaching a leaf (structure) in the recursions.
 - ▶ In each iteration, sample $0 \leq \alpha \leq 1$ uniformly at random.
 - ▶ Pop from the stack $top(i_R, j_R, i_S, j_S)$.
 - ▶ Pick a case of top according to α . For simplicity, we assume there is only one case here, i.e.

$$Q^{top} = \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1, k_2, j_S}^{left} Q_{k_1+1, j_R, i_S, k_2+1}^{right}$$

- ▶ Find k_1^*, k_2^* such that

$$\sum_{\substack{i_R \leq k_1 < k_1^* \\ i_S < k_2 \leq k_2^*}} \dots \simeq \alpha \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} \dots$$

- ▶ Push $left(i_R, k_1^*, k_2^*, j_S)$ and $right(k_1 + 1, j_R, i_S, k_2 + 1)$ onto the stack.



Sampling from the Boltzmann ensemble

- ▶ Push $l(1, n, 1, m)$ onto the stack.
- ▶ Iterate until the stack is empty, i.e. reaching a leaf (structure) in the recursions.
 - ▶ In each iteration, sample $0 \leq \alpha \leq 1$ uniformly at random.
 - ▶ Pop from the stack $top(i_R, j_R, i_S, j_S)$.
 - ▶ Pick a case of top according to α . For simplicity, we assume there is only one case here, i.e.

$$Q^{top} = \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1, k_2, j_S}^{left} Q_{k_1 + 1, j_R, i_S, k_2 + 1}^{right}$$

- ▶ Find k_1^*, k_2^* such that

$$\sum_{\substack{i_R \leq k_1 < k_1^* \\ i_S < k_2 \leq k_2^*}} \dots \simeq \alpha \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} \dots$$

- ▶ Push $left(i_R, k_1^*, k_2^*, j_S)$ and $right(k_1 + 1, j_R, i_S, k_2 + 1)$ onto the stack.



Sampling from the Boltzmann ensemble

- ▶ Push $l(1, n, 1, m)$ onto the stack.
- ▶ Iterate until the stack is empty, i.e. reaching a leaf (structure) in the recursions.
 - ▶ In each iteration, sample $0 \leq \alpha \leq 1$ uniformly at random.
 - ▶ Pop from the stack $top(i_R, j_R, i_S, j_S)$.
 - ▶ Pick a case of top according to α . For simplicity, we assume there is only one case here, i.e.

$$Q^{top} = \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1, k_2, j_S}^{left} Q_{k_1 + 1, j_R, i_S, k_2 + 1}^{right}$$

- ▶ Find k_1^*, k_2^* such that

$$\sum_{\substack{i_R \leq k_1 < k_1^* \\ i_S < k_2 \leq k_2^*}} \dots \simeq \alpha \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} \dots$$

- ▶ Push $left(i_R, k_1^*, k_2^*, j_S)$ and $right(k_1 + 1, j_R, i_S, k_2 + 1)$ onto the stack.



Sampling from the Boltzmann ensemble

- ▶ Push $l(1, n, 1, m)$ onto the stack.
- ▶ Iterate until the stack is empty, i.e. reaching a leaf (structure) in the recursions.
 - ▶ In each iteration, sample $0 \leq \alpha \leq 1$ uniformly at random.
 - ▶ Pop from the stack $top(i_R, j_R, i_S, j_S)$.
 - ▶ Pick a case of top according to α . For simplicity, we assume there is only one case here, i.e.

$$Q^{top} = \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1, k_2, j_S}^{left} Q_{k_1 + 1, j_R, i_S, k_2 + 1}^{right}$$

- ▶ Find k_1^*, k_2^* such that

$$\sum_{\substack{i_R \leq k_1 < k_1^* \\ i_S < k_2 \leq k_2^*}} \dots \simeq \alpha \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} \dots$$

- ▶ Push $left(i_R, k_1^*, k_2^*, j_S)$ and $right(k_1 + 1, j_R, i_S, k_2 + 1)$ onto the stack.



Sampling from the Boltzmann ensemble

- ▶ Push $l(1, n, 1, m)$ onto the stack.
- ▶ Iterate until the stack is empty, i.e. reaching a leaf (structure) in the recursions.
 - ▶ In each iteration, sample $0 \leq \alpha \leq 1$ uniformly at random.
 - ▶ Pop from the stack $top(i_R, j_R, i_S, j_S)$.
 - ▶ Pick a case of top according to α . For simplicity, we assume there is only one case here, i.e.

$$Q^{top} = \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1, k_2, j_S}^{left} Q_{k_1+1, j_R, i_S, k_2+1}^{right}$$

- ▶ Find k_1^*, k_2^* such that

$$\sum_{\substack{i_R \leq k_1 < k_1^* \\ i_S < k_2 \leq k_2^*}} \dots \simeq \alpha \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} \dots$$

- ▶ Push $left(i_R, k_1^*, k_2^*, j_S)$ and $right(k_1 + 1, j_R, i_S, k_2 + 1)$ onto the stack.



Sampling from the Boltzmann ensemble

- ▶ Push $l(1, n, 1, m)$ onto the stack.
- ▶ Iterate until the stack is empty, i.e. reaching a leaf (structure) in the recursions.
 - ▶ In each iteration, sample $0 \leq \alpha \leq 1$ uniformly at random.
 - ▶ Pop from the stack $top(i_R, j_R, i_S, j_S)$.
 - ▶ Pick a case of top according to α . For simplicity, we assume there is only one case here, i.e.

$$Q^{top} = \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1, k_2, j_S}^{left} Q_{k_1+1, j_R, i_S, k_2+1}^{right}$$

- ▶ Find k_1^*, k_2^* such that

$$\sum_{\substack{i_R \leq k_1 < k_1^* \\ i_S < k_2 \leq k_2^*}} \dots \simeq \alpha \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} \dots$$

- ▶ Push $left(i_R, k_1^*, k_2^*, j_S)$ and $right(k_1 + 1, j_R, i_S, k_2 + 1)$ onto the stack.



Sampling from the Boltzmann ensemble

- ▶ Push $l(1, n, 1, m)$ onto the stack.
- ▶ Iterate until the stack is empty, i.e. reaching a leaf (structure) in the recursions.
 - ▶ In each iteration, sample $0 \leq \alpha \leq 1$ uniformly at random.
 - ▶ Pop from the stack $top(i_R, j_R, i_S, j_S)$.
 - ▶ Pick a case of top according to α . For simplicity, we assume there is only one case here, i.e.

$$Q^{top} = \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} Q_{i_R, k_1, k_2, j_S}^{left} Q_{k_1+1, j_R, i_S, k_2+1}^{right}$$

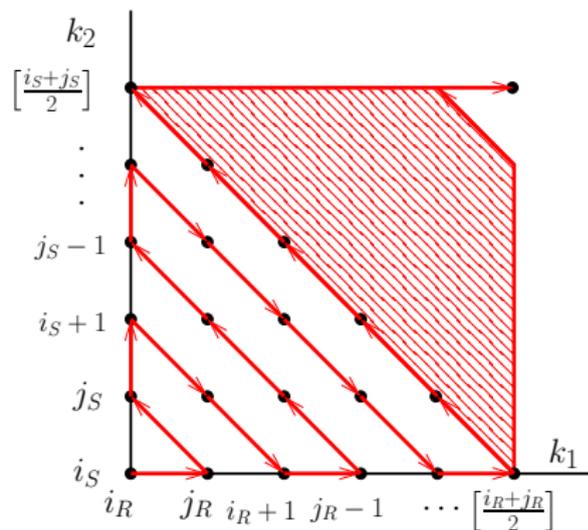
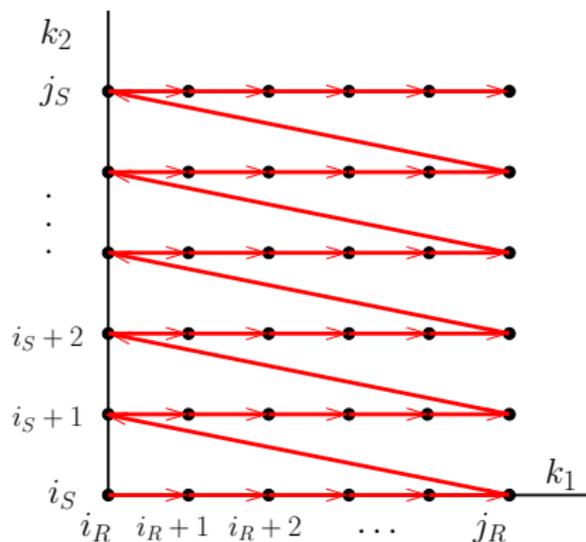
- ▶ Find k_1^*, k_2^* such that

$$\sum_{\substack{i_R \leq k_1 < k_1^* \\ i_S < k_2 \leq k_2^*}} \dots \simeq \alpha \sum_{\substack{i_R \leq k_1 < j_R \\ i_S < k_2 \leq j_S}} \dots$$

- ▶ Push $left(i_R, k_1^*, k_2^*, j_S)$ and $right(k_1 + 1, j_R, i_S, k_2 + 1)$ onto the stack.



Fast Ponty-style sampling of the Boltzmann ensemble



Time and space complexity of piRNA

- ▶ $O(n^4m^2 + n^2m^4)$ time.
- ▶ $O(n^2m^2)$ space.
- ▶ about 100 tables in the dynamic programming.
- ▶ takes about 1 day on 64 CPUs with 150GB RAM for two 110nt RNAs (OxyS-fhIA).

Therefore, a fast heuristic is on demand for high-throughput applications, possibly as a filtering step.



Time and space complexity of piRNA

- ▶ $O(n^4m^2 + n^2m^4)$ time.
- ▶ $O(n^2m^2)$ space.
- ▶ about 100 tables in the dynamic programming.
- ▶ takes about 1 day on 64 CPUs with 150GB RAM for two 110nt RNAs (OxyS-fhIA).

Therefore, a fast heuristic is on demand for high-throughput applications, possibly as a filtering step.



Binding sites prediction

biRNA : a fast algorithm to predict simultaneous binding sites of two nucleic acids

Pros

- ▶ Predicts multiple simultaneous binding sites.
- ▶ Computes a more accurate local energy of binding.
- ▶ Considers zigzags and crossing interactions.
- ▶ Maintains tractability for existing cases in the literature.

Cons

- ▶ Approximates the intramolecular site accessibility energy.
- ▶ Its running time grows exponentially with the maximum number of simultaneous binding sites.



Binding sites prediction

biRNA : a fast algorithm to predict simultaneous binding sites of two nucleic acids

Pros

- ▶ Predicts multiple simultaneous binding sites.
- ▶ Computes a more accurate local energy of binding.
- ▶ Considers zigzags and crossing interactions.
- ▶ Maintains tractability for existing cases in the literature.

Cons

- ▶ Approximates the intramolecular site accessibility energy.
- ▶ Its running time grows exponentially with the maximum number of simultaneous binding sites.



Steps of the algorithm for R and S

1. For all short subsequences W , compute $P_u(W)$, the prob. of being unpaired (Mückstein *et al.* 2008).
2. Obtain \mathcal{V} , a short list of candidate sites.
3. For all pairs W_1, W_2 , compute $P_u(W_1, W_2)$, the joint pairwise prob. of being simultaneously unpaired.
4. Build tree-structured Markov Random Fields (MRF) $\mathcal{T} = (\mathcal{V}, \mathcal{E})$ to approximate the distribution of being simultaneously unpaired (Chow and Liu 1968).
5. Compute $Q_{W^R W^S}^I$, the interaction partition functions restricted to subsequences W^R and W^S using piRNA.
6. Compute a matching between \mathcal{T}^R and \mathcal{T}^S that minimizes the binding energy or equivalently maximizes the binding probability.



Steps of the algorithm for R and S

1. For all short subsequences W , compute $P_u(W)$, the prob. of being unpaired (Mückstein *et al.* 2008).
2. Obtain \mathcal{V} , a short list of candidate sites.
3. For all pairs W_1, W_2 , compute $P_u(W_1, W_2)$, the joint pairwise prob. of being simultaneously unpaired.
4. Build tree-structured Markov Random Fields (MRF) $\mathcal{T} = (\mathcal{V}, \mathcal{E})$ to approximate the distribution of being simultaneously unpaired (Chow and Liu 1968).
5. Compute $Q_{W^R W^S}^I$, the interaction partition functions restricted to subsequences W^R and W^S using piRNA.
6. Compute a matching between \mathcal{T}^R and \mathcal{T}^S that minimizes the binding energy or equivalently maximizes the binding probability.



Steps of the algorithm for R and S

1. For all short subsequences W , compute $P_u(W)$, the prob. of being unpaired (Mückstein *et al.* 2008).
2. Obtain \mathcal{V} , a short list of candidate sites.
3. For all pairs W_1, W_2 , compute $P_u(W_1, W_2)$, the joint pairwise prob. of being simultaneously unpaired.
4. Build tree-structured Markov Random Fields (MRF) $\mathcal{T} = (\mathcal{V}, \mathcal{E})$ to approximate the distribution of being simultaneously unpaired (Chow and Liu 1968).
5. Compute Q_{WRWS}^I , the interaction partition functions restricted to subsequences W^R and W^S using piRNA.
6. Compute a matching between \mathcal{T}^R and \mathcal{T}^S that minimizes the binding energy or equivalently maximizes the binding probability.



Steps of the algorithm for R and S

1. For all short subsequences W , compute $P_u(W)$, the prob. of being unpaired (Mückstein *et al.* 2008).
2. Obtain \mathcal{V} , a short list of candidate sites.
3. For all pairs W_1, W_2 , compute $P_u(W_1, W_2)$, the joint pairwise prob. of being simultaneously unpaired.
4. Build tree-structured Markov Random Fields (MRF) $\mathcal{T} = (\mathcal{V}, \mathcal{E})$ to approximate the distribution of being simultaneously unpaired (Chow and Liu 1968).
5. Compute Q_{WRWS}^I , the interaction partition functions restricted to subsequences W^R and W^S using piRNA.
6. Compute a matching between \mathcal{T}^R and \mathcal{T}^S that minimizes the binding energy or equivalently maximizes the binding probability.



Steps of the algorithm for R and S

1. For all short subsequences W , compute $P_u(W)$, the prob. of being unpaired (Mückstein *et al.* 2008).
2. Obtain \mathcal{V} , a short list of candidate sites.
3. For all pairs W_1, W_2 , compute $P_u(W_1, W_2)$, the joint pairwise prob. of being simultaneously unpaired.
4. Build tree-structured Markov Random Fields (MRF) $\mathcal{T} = (\mathcal{V}, \mathcal{E})$ to approximate the distribution of being simultaneously unpaired (Chow and Liu 1968).
5. Compute $Q_{W^R W^S}^I$, the interaction partition functions restricted to subsequences W^R and W^S using piRNA.
6. Compute a matching between \mathcal{T}^R and \mathcal{T}^S that minimizes the binding energy or equivalently maximizes the binding probability.



Steps of the algorithm for R and S

1. For all short subsequences W , compute $P_u(W)$, the prob. of being unpaired (Mückstein *et al.* 2008).
2. Obtain \mathcal{V} , a short list of candidate sites.
3. For all pairs W_1, W_2 , compute $P_u(W_1, W_2)$, the joint pairwise prob. of being simultaneously unpaired.
4. Build tree-structured Markov Random Fields (MRF) $\mathcal{T} = (\mathcal{V}, \mathcal{E})$ to approximate the distribution of being simultaneously unpaired (Chow and Liu 1968).
5. Compute $Q_{W^R W^S}^I$, the interaction partition functions restricted to subsequences W^R and W^S using `pirNA`.
6. Compute a matching between \mathcal{T}^R and \mathcal{T}^S that minimizes the binding energy or equivalently maximizes the binding probability.



biRNA

Binding energy minimization

Exhaustive search to find matching

$M = \{(W_1^R, W_1^S), (W_2^R, W_2^S), \dots, (W_k^R, W_k^S)\}$ that minimizes

$$\Delta G(M) = ED_u^R(M) + ED_u^S(M) + \Delta G_b^{RS}(M),$$

in which

$$ED_u^R(M) = -RT \log P_u^{R*}(W_1^R, W_2^R, \dots, W_k^R)$$

$$ED_u^S(M) = -RT \log P_u^{S*}(W_1^S, W_2^S, \dots, W_k^S)$$

$$\Delta G_b^{RS}(M) = -RT \sum_{1 \leq i \leq k} \log(Q_{W_i^R W_i^S}^I - Q_{W_i^R} Q_{W_i^S}).$$

R is the universal gas constant and T is temperature.



Experimental results

Multi-sites

Pair	Binding Site(s)		biRNA		RNAup	
	Literature		Site(s)		Site	
OxyS-fhlA	[22,30]	[95,87]	(23,30)	(94,87)	-	-
	[98,104]	[45,39]	(96,104)	(48,39)	(96,104)	(48,39)
CopA-CopT	[22,33]	[70,59]	(22,31)	(70,61)	-	-
	[48,56]	[44,36]	(49,57)	(43,35)	(49,67)	(43,24)
	[62,67]	[29,24]	(58,67)	(33,24)	-	-



Experimental results

Uni-sites

Pair	
GcvB	gltI
GcvB	argT
GcvB	dppA
GcvB	livJ
GcvB	livK
GcvB	oppA
GcvB	STM4351
MicA	lamB
MicA	ompA
DsrA	rpoS
RprA	rpoS
IstR	tisA
MicC	ompC
MicF	ompF
RyhB	sdhD
RyhB	sodB
SgrS	ptsG
IncRNA ₅₄	repZ

Lengths: 71-253 nt

Running time: 10 min - 1 hour on 8 dual core CPUs and 20GB of RAM



Summary

- ▶ We presented piRNA an $O(n^4 m^2 + n^2 m^4)$ -time $O(n^2 m^2)$ -space complexity algorithm for interaction partition function, base-pair probabilities, minimum free energy secondary structure, equilibrium concentrations, melting temperature, and some other derivatives of the partition function.
- ▶ piRNA outperforms all other alternatives and is available at <http://compbio.cs.wayne.edu/chitsaz/>.
- ▶ We presented biRNA , a fast RNA-RNA binding sites prediction algorithm.
- ▶ biRNA 's tree-structured MRF approximation is accurate enough for predicting binding sites and may be used in other applications.



Summary

- ▶ We presented piRNA an $O(n^4 m^2 + n^2 m^4)$ -time $O(n^2 m^2)$ -space complexity algorithm for interaction partition function, base-pair probabilities, minimum free energy secondary structure, equilibrium concentrations, melting temperature, and some other derivatives of the partition function.
- ▶ piRNA outperforms all other alternatives and is available at <http://compbio.cs.wayne.edu/chitsaz/>.
- ▶ We presented biRNA , a fast RNA-RNA binding sites prediction algorithm.
- ▶ biRNA 's tree-structured MRF approximation is accurate enough for predicting binding sites and may be used in other applications.



Summary

- ▶ We presented piRNA an $O(n^4 m^2 + n^2 m^4)$ -time $O(n^2 m^2)$ -space complexity algorithm for interaction partition function, base-pair probabilities, minimum free energy secondary structure, equilibrium concentrations, melting temperature, and some other derivatives of the partition function.
- ▶ piRNA outperforms all other alternatives and is available at <http://compbio.cs.wayne.edu/chitsaz/>.
- ▶ We presented biRNA , a fast RNA-RNA binding sites prediction algorithm.
- ▶ biRNA 's tree-structured MRF approximation is accurate enough for predicting binding sites and may be used in other applications.



Summary

- ▶ We presented piRNA an $O(n^4 m^2 + n^2 m^4)$ -time $O(n^2 m^2)$ -space complexity algorithm for interaction partition function, base-pair probabilities, minimum free energy secondary structure, equilibrium concentrations, melting temperature, and some other derivatives of the partition function.
- ▶ piRNA outperforms all other alternatives and is available at <http://compbio.cs.wayne.edu/chitsaz/>.
- ▶ We presented biRNA , a fast RNA-RNA binding sites prediction algorithm.
- ▶ biRNA 's tree-structured MRF approximation is accurate enough for predicting binding sites and may be used in other applications.



Future work

- ▶ RNA design for positive and negative interactions.
- ▶ Better interaction energy model, which requires more data.
- ▶ Incorporation of non-canonical base pairs.



Future work

- ▶ RNA design for positive and negative interactions.
- ▶ Better interaction energy model, which requires more data.
- ▶ Incorporation of non-canonical base pairs.



Future work

- ▶ RNA design for positive and negative interactions.
- ▶ Better interaction energy model, which requires more data.
- ▶ Incorporation of non-canonical base pairs.



Acknowledgement

Collaborators

- ▶ Rolf Backofen, University of Freiburg, Germany
- ▶ Cenk Sahinalp, SFU, Canada
- ▶ Raheleh Salari

Funding



Michael Smith Foundation for
Health Research



MITACS

BCID
Bioinformatics for Combating Infectious Diseases

Deutsche
Forschungsgemeinschaft

DFG

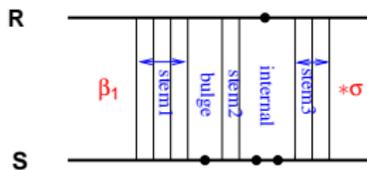


Thanks for your attention!



Hybrid component

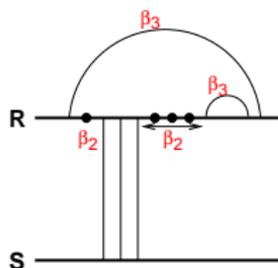
Example



$$G^{\text{hybrid}} = \beta_1 + \sigma(G^{\text{stem}_1} + G^{\text{bulge}} + G^{\text{stem}_2} + G^{\text{internal}} + G^{\text{stem}_3}).$$

Kissing loop

Example



$$G^{\text{kissing}} = 4\beta_2 + 2\beta_3.$$