

Efficient Computation of Base-Pairing Probabilities in Multi-Strand RNA Folding

Ronny Lorenz
ronny@tbi.univie.ac.at

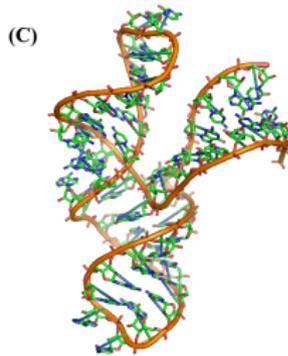
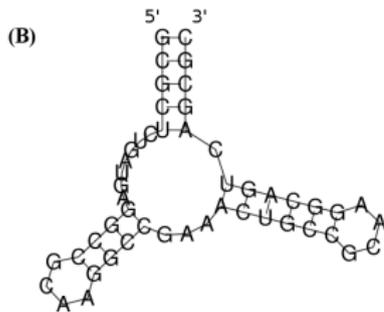
University of Vienna, Theoretical Biochemistry Group (TBI)

Valletta, Malta, February 24th 2020

RNA Secondary Structures

RNAs fold hierarchical (A) → (B) → (C)

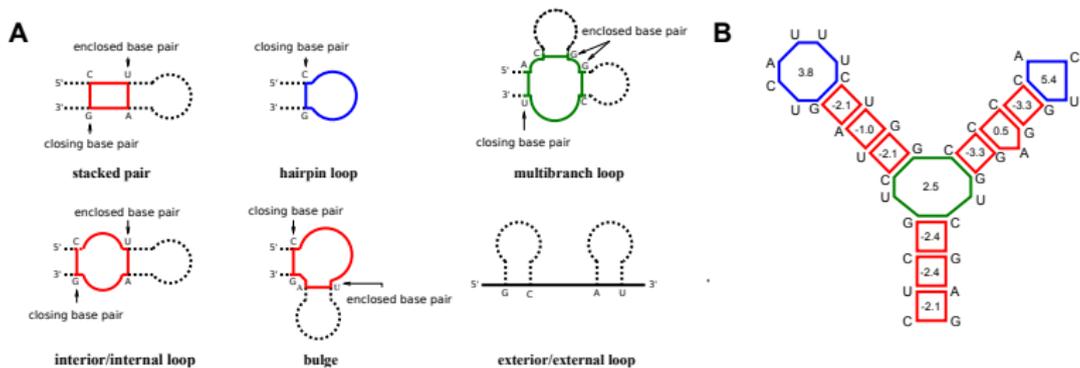
(A) 5' - GCGCUCUGAUGAGGCCGCAAGGCCGAAACUGCCGCAAGGCAGUCAGCGC - 3'



Secondary Structure (B):

- Set of nested base pairs
- Captures majority of stabilizing interactions
- Many thermodynamic properties can be predicted efficiently
- Very good prediction accuracy for small RNAs
- Accuracy drops to 40%-70% for longer sequences

RNA Secondary Structures - Loops vs. Base Pairs



- secondary structures s can be uniquely decomposed into loops L
- stabilizing energy contributions (mostly) from stacked base pairs
- destabilizing contributions from unpaired bases in loops
- each loop L is assigned a free energy contribution E_L ¹

$$E(s) \approx \sum_{L \in s} E_L$$

¹Turner, DH & Mathews, DH (2009). NNDB: The nearest neighbor parameter database for predicting stability of nucleic acid secondary structure., Nucleic Acids Research 38, D280-D282

RNA Secondary Structures - Statistical Thermodynamics

$$p(F) \propto e^{-\frac{E(F)}{RT}}$$

Most probable structure:

$$MFE = \min_s E(s)$$

Partition Function:

$$Q = \sum_s e^{-\frac{E(s)}{RT}}$$

Probability of a structure:

$$p(s) = \frac{e^{-E(s)/RT}}{Q}$$

Probability of Base Pair (k, l) :

$$p_{k,l} = \frac{1}{Q} \sum_{s|(k,l) \in s} e^{-\frac{E(s)}{RT}}$$

RNA Secondary Structures - Prediction

Dynamic Programming (DP) algorithm²



²Nussinov, R & Pieczenik, G & Griggs, JR and Kleitman, DJ (1978). Algorithms for Loop Matchings., SIAM J. Appl. Math., 35(1), 68-82

RNA Secondary Structures - Prediction

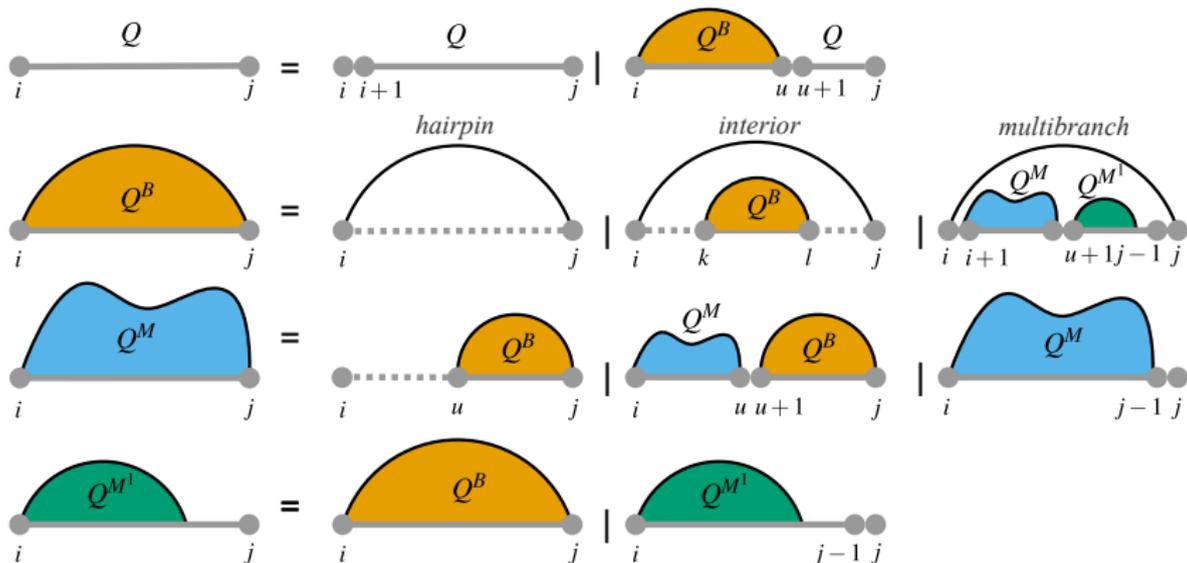
Dynamic Programming (DP) algorithm²



²Nussinov, R & Pieczenik, G & Griggs, JR and Kleitman, DJ (1978). Algorithms for Loop Matchings., SIAM J. Appl. Math., 35(1), 68-82

RNA Secondary Structures - Prediction

Dynamic Programming (DP) algorithm^{2 3}

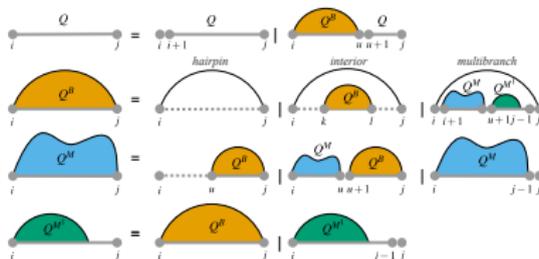


²Zuker M and Stiegler P (1981). Optimal computer folding of large RNA sequences using thermodynamics and auxiliary information., *Nucleic Acids Res*, 9(1):133-148

³McCaskill JS (1990). The equilibrium partition function and base pair binding probabilities for RNA secondary structure., *Biopolymers*, 29(6-7):1105-1119

RNA Secondary Structures - Prediction

Dynamic Programming (DP) algorithm



$$Q_{i,j} = Q_{i+1,j} + \sum_{i < u < j} Q_{i,u}^B Q_{u+1,j}$$

$$Q_{i,j}^B = e^{-\frac{H(i,j)}{RT}} + \sum_{i < k < l < j} e^{-\frac{I(i,j,k,l)}{RT}} Q_{k,l}^B + e^{-\frac{a+b}{RT}} \sum_{i < u < j} Q_{i+1,u}^M Q_{u+1,j-1}^{M^1}$$

$$Q_{i,j}^M = \sum_{i \leq u < j} e^{-\frac{(u-i)c+b}{RT}} Q_{u,j}^B + e^{-\frac{b}{RT}} \sum_{i < u < j} Q_{i,u}^M Q_{u+1,j}^B + e^{-\frac{c}{RT}} Q_{i,j-1}^{M^1}$$

$$Q_{i,j}^{M^1} = e^{-\frac{b}{RT}} Q_{i,j}^B + e^{-\frac{c}{RT}} Q_{i,j-1}^{M^1}$$

Asymptotic complexity: $O(n^3)$ time and $O(n^2)$ memory

Multiple Interacting Nucleic Acid Strands²

Straight-forward extension of single sequence case

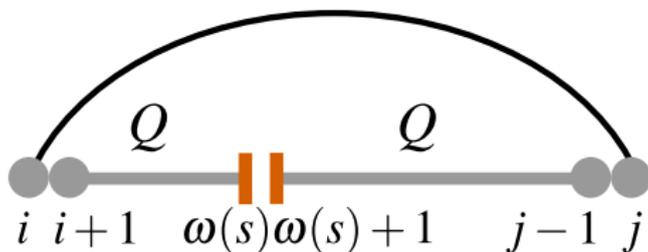
- consider complexes of N strands, i.e. one connected component
- restrict state space to intermolecular base pairs w/o crossings
- concatenate all strands ($n = n_1 + n_2 + \dots + n_N$)
- prohibit strand nicks in “regular” loops
- treat cases with nicks as “external” loops w/ additional rule
- process all non-cyclic permutations π of strand concatenations
- correct for overcounting of symmetric cases

²Dirks, RM, Bois, JS, Schaeffer, JM, Winfree, E, and Pierce, NA (2007).
Thermodynamic analysis of interacting nucleic acid strands., SIAM Rev., 49:65–88.

Multiple Interacting Nucleic Acid Strands²

Straight-forward extension of single sequence case

- consider complexes of N strands, i.e. one connected component
- restrict state space to intermolecular base pairs w/o crossings
- concatenate all strands ($n = n_1 + n_2 + \dots + n_N$)
- prohibit strand nicks in “regular” loops
- treat cases with nicks as “external” loops w/ **additional rule**
- process all non-cyclic permutations π of strand concatenations
- correct for overcounting of symmetric cases

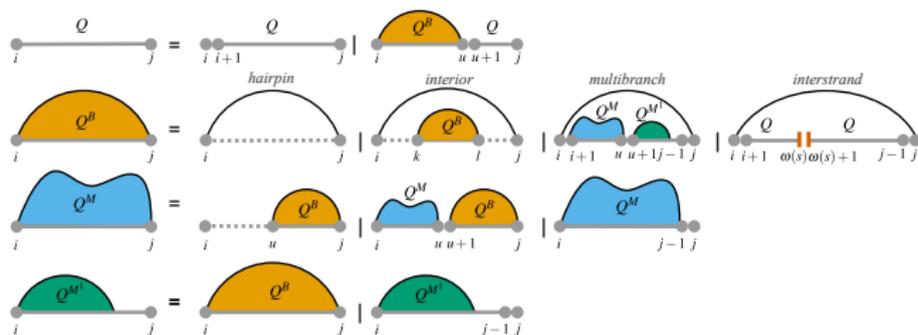


²Dirks, RM, Bois, JS, Schaeffer, JM, Winfree, E, and Pierce, NA (2007).
Thermodynamic analysis of interacting nucleic acid strands., SIAM Rev., 49:65–88.

Multiple Interacting Nucleic Acid Strands²

Straight-forward extension of single sequence case

- consider complexes of N strands, i.e. one connected component
- restrict state space to intermolecular base pairs w/o crossings
- concatenate all strands ($n = n_1 + n_2 + \dots + n_N$)
- prohibit strand nicks in “regular” loops
- treat cases with nicks as “external” loops w/ **additional rule**
- process all non-cyclic permutations π of strand concatenations
- correct for overcounting of symmetric cases



²Dirks, RM, Bois, JS, Schaeffer, JM, Winfree, E, and Pierce, NA (2007). Thermodynamic analysis of interacting nucleic acid strands., SIAM Rev., 49:65–88.

Interacting Nucleic Acid Strands - Base Pairing Probabilities

For *single* strand case (DP algorithm with $O(n^3)$ time, $O(n^2)$ memory):

$$\begin{aligned} p_{k,l} &= \frac{1}{Q} \sum_{s|(k,l) \in s} e^{-\frac{E(s)}{RT}} \\ &= \frac{1}{Q} Q_{k,l}^B \hat{Q}_{k,l}, \quad \text{with} \\ \hat{Q}_{k,l} &= \underbrace{\bar{Q}_{k,l}}_{\text{not enclosed by any bp}} + \underbrace{\check{Q}_{k,l}}_{\text{enclosed by bp}} \end{aligned}$$

Interacting Nucleic Acid Strands - Base Pairing Probabilities

For *single* strand case (DP algorithm with $O(n^3)$ time, $O(n^2)$ memory):

$$\begin{aligned}
 p_{k,l} &= \frac{1}{Q} \sum_{s|(k,l) \in s} e^{-\frac{E(s)}{RT}} \\
 &= \frac{1}{Q} Q_{k,l}^B \widehat{Q}_{k,l}, \quad \text{with} \\
 \widehat{Q}_{k,l} &= \underbrace{\bar{Q}_{k,l}}_{\text{not enclosed by any bp}} + \underbrace{\check{Q}_{k,l}}_{\text{enclosed by bp}}
 \end{aligned}$$

For complexes of N strands:

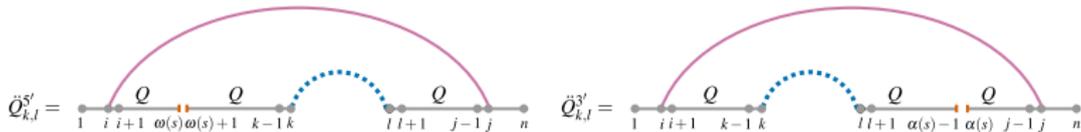
$$\begin{aligned}
 p_{k,l} &= \sum_{\pi} w(\pi) p_{k,l}[\pi] \\
 &= \frac{1}{Q} \sum_{\pi} \widehat{Q}_{k,l}[\pi] Q_{k,l}^B[\pi], \quad \text{with} \\
 \widehat{Q}_{k,l}[\pi] &= \underbrace{\bar{Q}_{k,l}[\pi]}_{\text{not enclosed by any bp}} + \underbrace{\check{Q}_{k,l}[\pi]}_{\text{enclosed by bp}} + \underbrace{\ddot{Q}_{k,l}[\pi]}_{\text{enclosed by bp w/ nick in loop}}
 \end{aligned}$$

What is the asymptotic complexity to compute $\ddot{Q}_{k,l}[\pi]$?

Interacting Nucleic Acid Strands - Base Pairing Probabilities

Additional case for nicked loops (for particular π):

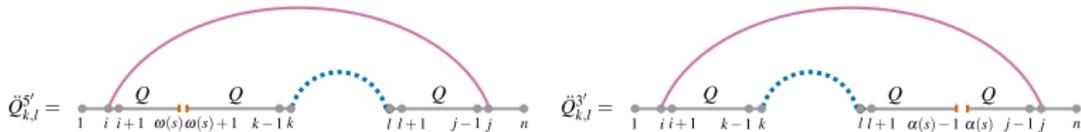
$$\ddot{Q}_{k,l} = \ddot{Q}_{k,l}^{5'} + \ddot{Q}_{k,l}^{3'}$$



Interacting Nucleic Acid Strands - Base Pairing Probabilities

Additional case for nicked loops (for particular π):

$$\ddot{Q}_{k,l} = \ddot{Q}_{k,l}^{5'} + \ddot{Q}_{k,l}^{3'}$$

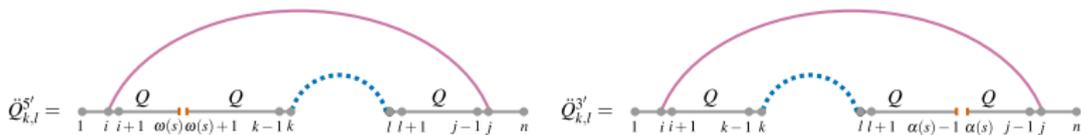


Computing all $\ddot{Q}_{k,l}$ by considering all enclosing pairs (i, j) and $N-1$ strand nicks seems to require $O(n^4 N)$ operations

Interacting Nucleic Acid Strands - Base Pairing Probabilities

Additional case for nicked loops (for particular π):

$$\ddot{Q}_{k,l} = \ddot{Q}_{k,l}^{5'} + \ddot{Q}_{k,l}^{3'}$$



Computing all $\ddot{Q}_{k,l}$ by considering all enclosing pairs (i,j) and $N-1$ strand nicks seems to require $O(n^4N)$ operations

- Dirks et al., 2007:

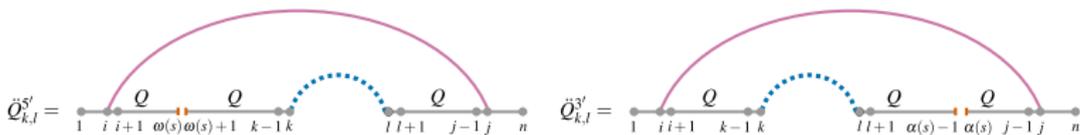
“... equilibrium probability of each intrastrand and interstrand base pair ... can be calculated by backtracking through the partition function algorithm ... applying a particular algorithmic transformation at each step”
- Wolfe et al., 2017:

“... the equilibrium base-pairing properties ... must be calculated for each complex $j \in \Psi$ using $\Theta(|\phi_j|^3)$ dynamic programs.” (here, $|\phi_j| \equiv n$)

Interacting Nucleic Acid Strands - Base Pairing Probabilities

Additional case for nicked loops (for particular π):

$$\ddot{Q}_{k,l} = \ddot{Q}_{k,l}^{5'} + \ddot{Q}_{k,l}^{3'}$$



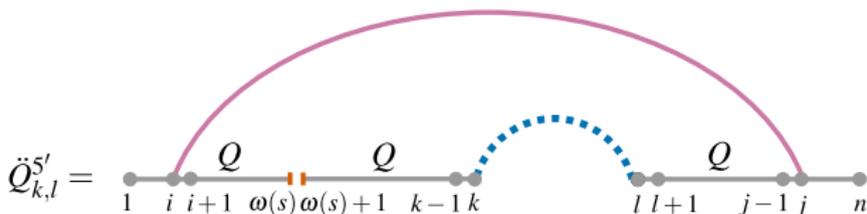
Computing all $\ddot{Q}_{k,l}$ by considering all enclosing pairs (i,j) and $N-1$ strand nicks seems to require $O(n^4N)$ operations

- Dirks et al., 2007:
“... equilibrium probability of each intrastrand and interstrand base pair ... can be calculated by backtracking through the partition function algorithm ... applying a particular algorithmic transformation at each step”
- Wolfe et al., 2017:
“... the equilibrium base-pairing properties ... must be calculated for each complex $j \in \Psi$ using $\Theta(|\phi_j|^3)$ dynamic programs.” (here, $|\phi_j| \equiv n$)

Still, no reference to the algorithm, so how to achieve that?

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



$$\ddot{Q}_{k,l}^{5'} = \sum_{\substack{1 \leq i < k \\ l < j \leq n}} \hat{Q}_{i,j} Q_{l+1,j-1} \times \sum_{s | i \leq \omega(s) < k} Q_{i+1,\omega(s)} Q_{\omega(s)+1,k-1}$$

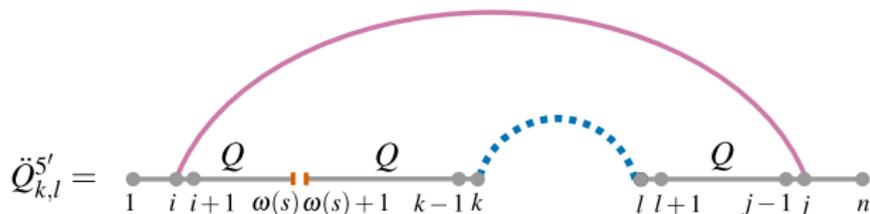
Apply Dynamic Programming paradigm:

- Trade computation time against memory consumption³
- Extract parts that are computed redundantly for different $\ddot{Q}_{k,l}^{5'}$

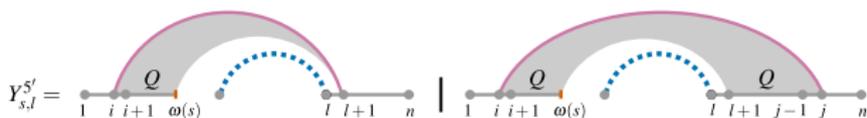
³Similar to McCaskill, JS (1990). The equilibrium partition function and base pair binding probabilities for RNA secondary structure., Biopolymers, 29:1105-1119.

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



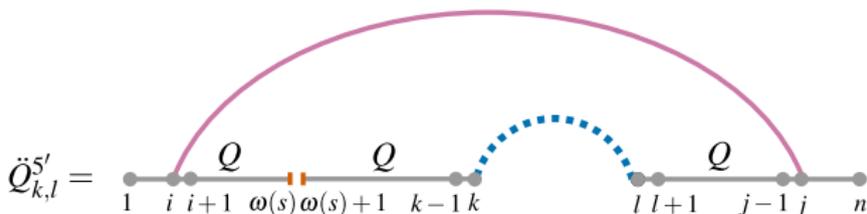
1st step (fixed l):



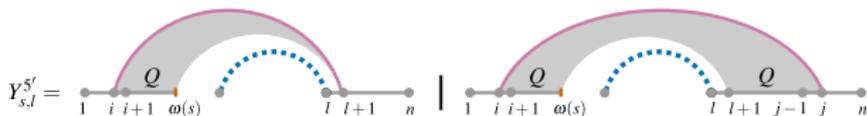
- pre-compute “enclosed” part $Y_{s,l}^{5'}$ up to $\omega(s)$
- re-use $Y_{s,l}^{5'}$ for all k

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



1st step (fixed l):

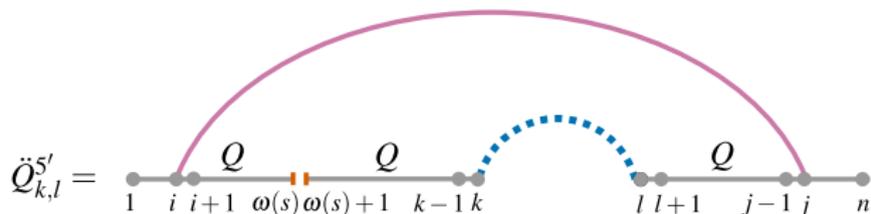


$$\ddot{Q}_{k,l}^{5'} = Y_{\sigma(k-1),l}^{5'} + \sum_{s|\omega(s)<k} Q_{\omega(s)+1,k-1} Y_{s,l}^{5'} \quad (\text{indep. of } i \text{ and } j)$$

$$Y_{s,l}^{5'} = \sum_{j>l} Q_{l+1,j-1} \times \left(\hat{Q}_{\omega(s),j} + \sum_{i<\omega(s)} \hat{Q}_{i,j} Q_{i+1,\omega(s)} \right) \quad (\text{indep. of } k)$$

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



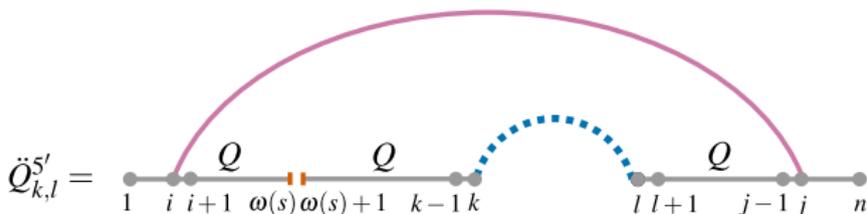
1st step (fixed l):



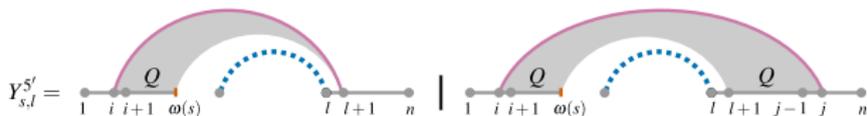
Can we do better than $O(n^3N)$?

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



1st step (fixed l):



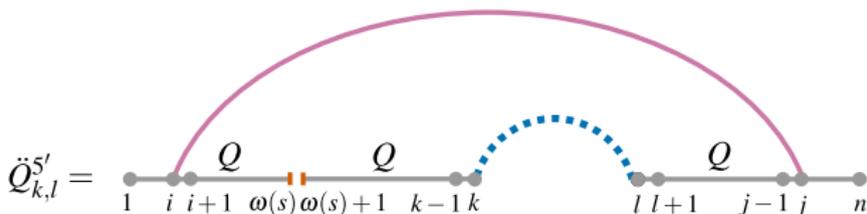
Can we do better than $O(n^3N)$?

Observations when comparing $Y_{s,l-1}^{5'}$ against $Y_{s,l}^{5'}$:

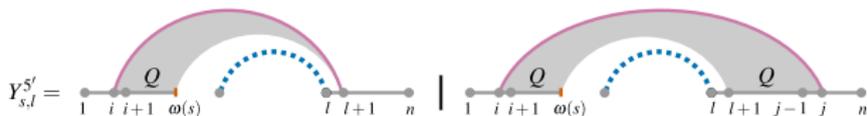
- “left” contribution stays the same
- “right” contribution includes $Q_{l,j}$ instead of $Q_{l+1,j}$
- one more j to account for ($j = l$)

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



1st step (fixed l):



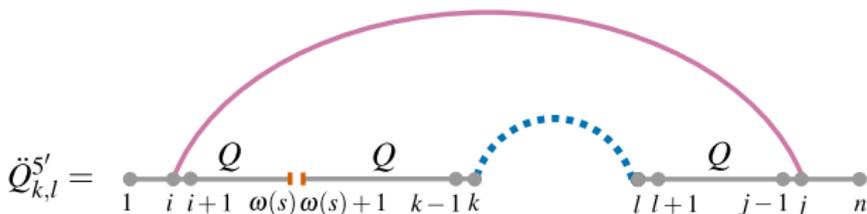
Can we do better than $O(n^3N)$?

Observations when comparing $Y_{s,l-1}^{5'}$ against $Y_{s,l}^{5'}$:

- “left” contribution stays the same (**pre-compute and re-use!**)
- “right” contribution includes $Q_{l,j}$ instead of $Q_{l+1,j}$
- one more j to account for ($j = l$)

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



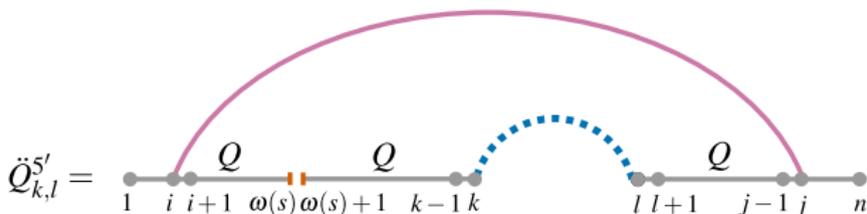
2nd step (pre-compute “left” part of $Y_{s,l}^{5'}$):



- “left” part ($Y_{s,j}^{5''}$) delimited by $\omega(s)$ and j is independent of k and l
- re-use $Y_{s,j}^{5''}$ to compute $Y_{s,l}^{5'}$ for all l

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



2nd step (pre-compute “left” part of $Y_{s,l}^{5'}$):

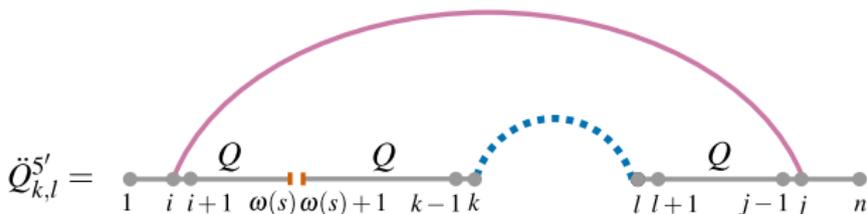


$$Y_{s,l}^{5'} = Y_{s,l+1}^{5''} + \sum_{j>l+1} Q_{l+1,j-1} \cdot Y_{s,j}^{5''} \quad (\text{indep. of } i \text{ and } k)$$

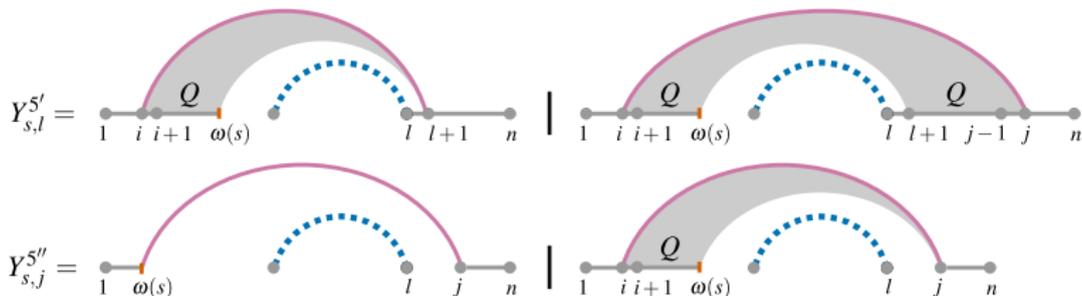
$$Y_{s,j}^{5''} = \hat{Q}_{\omega(s),j} + \sum_{i<\omega(s)} \hat{Q}_{i,j} \cdot Q_{i+1,\omega(s)} \quad (\text{indep. of } k \text{ and } l)$$

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



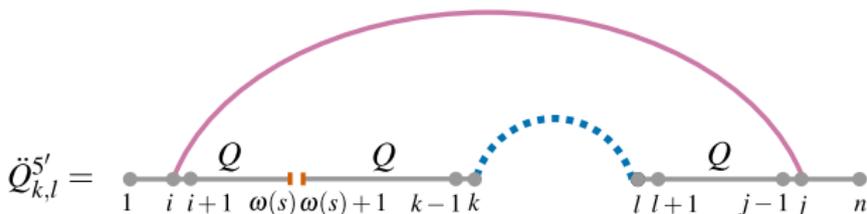
Finally:



Complexity: $O(n^2N)$ time and additional $O(nN)$ memory

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 0: strand nick is on 5' side



Finally:

$$\ddot{Q}_{k,l}^{5'} = Y_{\sigma(k-1),l}^{5'} + \sum_{s|\omega(s)<k} Q_{\omega(s)+1,k-1} Y_{s,l}^{5'}$$

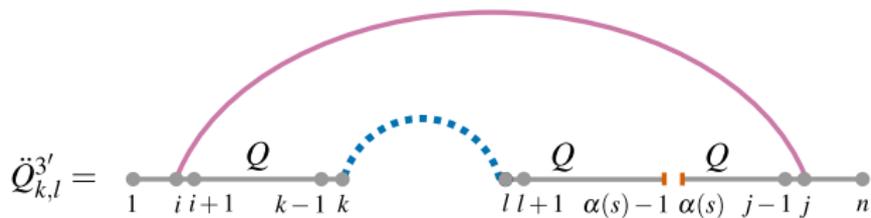
$$Y_{s,l}^{5'} = Y_{s,l+1}^{5''} + \sum_{j>l+1} Q_{l+1,j-1} \cdot Y_{s,j}^{5''}$$

$$Y_{s,j}^{5''} = \hat{Q}_{\omega(s),j} + \sum_{i<\omega(s)} \hat{Q}_{i,j} \cdot Q_{i+1,\omega(s)}$$

Complexity: $O(n^2N)$ time and additional $O(nN)$ memory

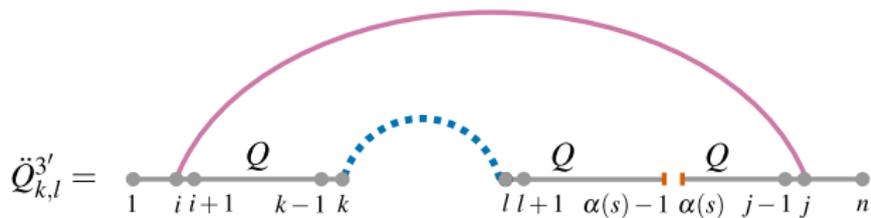
Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 1: strand nick is on 3' side



Interacting Nucleic Acid Strands - Base Pairing Probabilities

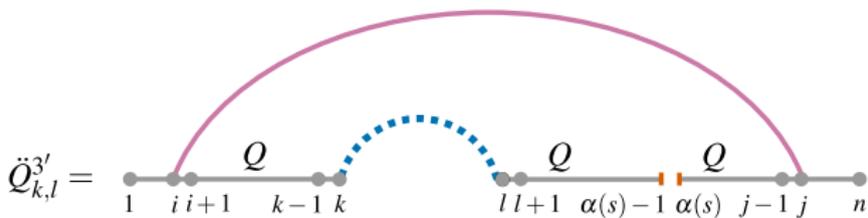
case 1: strand nick is on 3' side



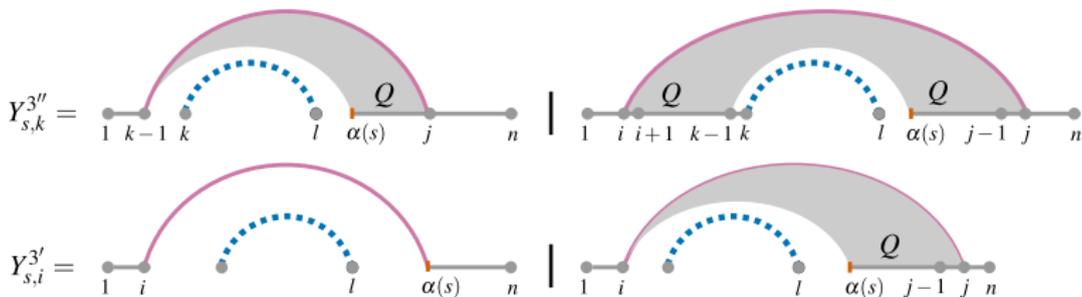
Using similar algorithmic transformations as for $\ddot{Q}_{k,l}^{5'}$

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 1: strand nick is on 3' side

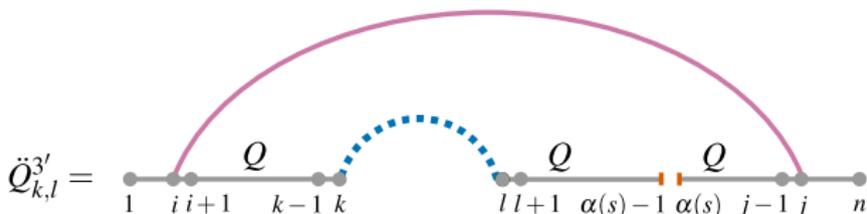


Using similar algorithmic transformations as for $\ddot{Q}_{k,l}^{5'}$



Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 1: strand nick is on 3' side



Using similar algorithmic transformations as for $\hat{Q}_{k,l}^{5'}$

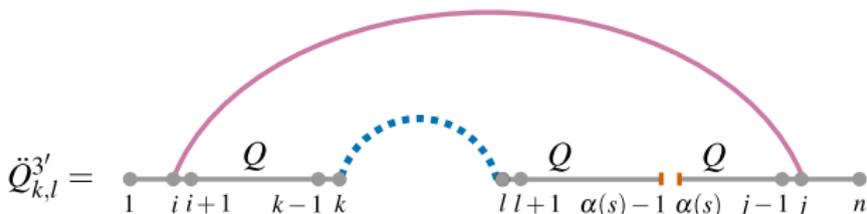
$$\hat{Q}_{k,l}^{3'} = Y_{\sigma(l+1),k}^{3''} + \sum_{s|\alpha(s)>l+1} Q_{l+1,\alpha(s)-1} Y_{s,k}^{3''}$$

$$Y_{s,k}^{3''} = \sum_{i<k} Q_{i+1,k-1} Y_{s,i}^{3'}$$

$$Y_{s,i}^{3'} = \hat{Q}_{i,\alpha(s)} + \sum_{j>\alpha(s)} \hat{Q}_{i,j} Q_{\alpha(s),j-1}$$

Interacting Nucleic Acid Strands - Base Pairing Probabilities

case 1: strand nick is on 3' side



Using similar algorithmic transformations as for $\hat{Q}_{k,l}^{5'}$

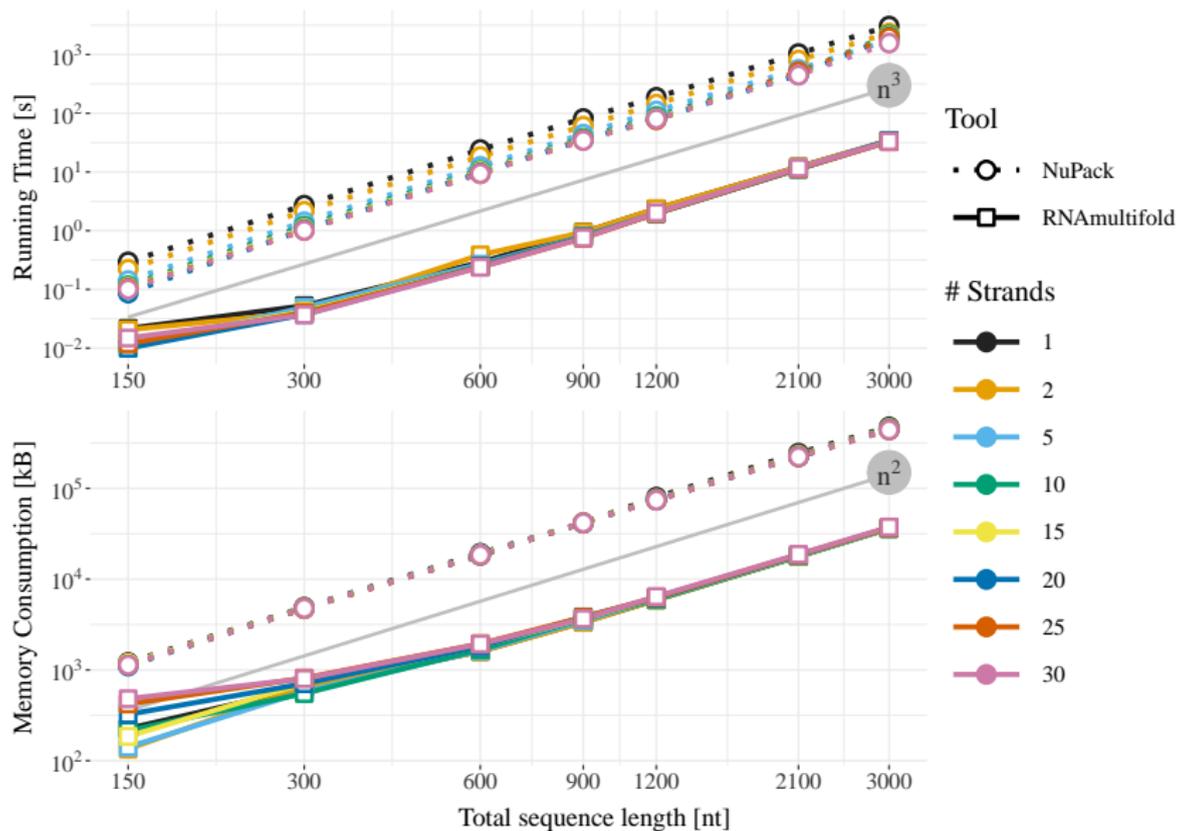
$$\hat{Q}_{k,l}^{3'} = Y_{\sigma(l+1),k}^{3''} + \sum_{s|\alpha(s)>l+1} Q_{l+1,\alpha(s)-1} Y_{s,k}^{3''}$$

$$Y_{s,k}^{3''} = \sum_{i<k} Q_{i+1,k-1} Y_{s,i}^{3'}$$

$$Y_{s,i}^{3'} = \hat{Q}_{i,\alpha(s)} + \sum_{j>\alpha(s)} \hat{Q}_{i,j} Q_{\alpha(s),j-1}$$

Total effort for all $\hat{Q}_{k,l}$: $O(n^2N)$ time using additional $O(nN)$ memory

Runtime and Memory Consumption



Conclusion

- overhead of “nicked” loops is negligible
- all $p_{k,l}$ can indeed be computed in $O(|\Pi|n^3)$ time
- implementation available as `RNAmultifold`⁴
- 50 – 65× faster, 7× less memory than NUPACK 3.2.2 (per π)
- full constraints support, e.g. to restrict state space or include experimental probing data (SHAPE, etc.)
- intramolecular G-Quadruplex

Outlook

- automatically compute over all π and complexes of size N
- add MFE, Boltzmann sampling, and suboptimal enumeration
- include ligand binding support, e.g. SSB proteins
- add concentration dependency
- re-use (parts of) DP tables for different π

⁴ViennaRNA Package 2.5.0alpha

Thanks to

- **Christoph Flamm**
- **Ivo Hofacker**
- **Peter Stadler**

Thank You for your attention!

This work was funded in parts by the German Federal Ministry of Education and Research (BMBF, project no. 031A538A, de.NBI-RBC, to PFS and project no. 031L0164C, RNAProNet), and the Austrian science fund FWF (project no. I2674 "Prediction of RNA-RNA interactions", project no. F 43 "RNA regulation of the transcriptome").