

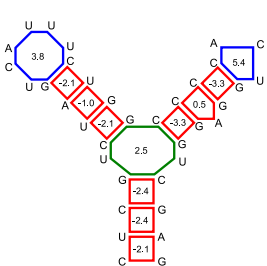
Modified Nucleotides in RNA structure prediction

Ronny Lorenz

University of Leipzig, Computer Science (Chair for bioinformatics)
University of Vienna, Theoretical Biochemistry Group (TBI)

Benasque, Spain, July 23, 2024

RNA Secondary Structures - Nearest Neighbor Energy Model



$$\begin{aligned}
 E(\gamma) &= \begin{matrix} \text{U} & \text{A} \\ \text{C} & \text{G} \end{matrix} \boxed{-2.1} + \begin{matrix} \text{C} & \text{G} \\ \text{U} & \text{A} \end{matrix} \boxed{-2.4} + \begin{matrix} \text{G} & \text{C} \\ \text{C} & \text{G} \end{matrix} \boxed{-2.4} + \begin{matrix} \text{G} & \text{C} \\ \text{U} & \text{G} \end{matrix} \boxed{2.5} + \\
 &\begin{matrix} \text{G} & \text{U} \\ \text{C} & \text{G} \end{matrix} \boxed{-2.1} + \begin{matrix} \text{U} & \text{G} \\ \text{A} & \text{U} \end{matrix} \boxed{-1.0} + \begin{matrix} \text{C} & \text{U} \\ \text{G} & \text{A} \end{matrix} \boxed{-2.1} + \begin{matrix} \text{A} & \text{U} \\ \text{C} & \text{U} \end{matrix} \boxed{3.8} + \\
 &\begin{matrix} \text{C} & \text{G} \\ \text{C} & \text{G} \end{matrix} \boxed{-3.3} + \begin{matrix} \text{C} & \text{G} \\ \text{G} & \text{A} \end{matrix} \boxed{0.5} + \begin{matrix} \text{C} & \text{G} \\ \text{G} & \text{G} \end{matrix} \boxed{-3.3} + \begin{matrix} \text{A} & \text{C} \\ \text{C} & \text{U} \end{matrix} \boxed{5.4} \\
 &= -6.50 \text{ kcal/mol}
 \end{aligned}$$

- Secondary structures s can be uniquely decomposed into loops L
- Contributions of a base pair only depends on neighboring pairs
- Each loop L is assigned a free energy contribution E_L ^[1]

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

[1] A. Mittal et al. "NNDB: An Expanded Database of Nearest Neighbor Parameters for Predicting Stability of Nucleic Acid Secondary Structures". In: *Journal of Molecular Biology* (2024), p. 168549. DOI: 10. 1016/j. jmb. 2024. 168549

RNA Secondary Structures - Statistical Thermodynamics

$$\begin{aligned} p(F) &\propto e^{-\beta E(F)}, \quad \text{with} \quad \beta = \frac{1}{RT} \\ &= \frac{1}{Q} \sum_{s|F \in s} e^{-\beta E(s)}, \quad \text{with} \quad Q = \sum_s e^{-\beta E(s)} \end{aligned}$$

Efficient dynamic programming algorithms are available to compute

- Minimum free energy $MFE = \min_s E(s)$
- Partition function Q
- Probabilities p_{kl} for base pairs (k, l) , i.e. $p_{kl} = p(F)$ with $F = (k, l)$
- Suboptimal and locally stable structures, ligand binding, ...

Implementations for different input data types^[2]

- single sequences
- (multiple) interacting sequences
- sequence alignments (consensus structure)

[2] R. Lorenz et al. "ViennaRNA Package 2.0". In: *Algorithms for Molecular Biology* 6 (2011), pp. 1–14. DOI: 10. 1186/ 1748-7188-6-26

RNA Secondary Structures - Statistical Thermodynamics

$$\begin{aligned} p(F) &\propto e^{-\beta E(F)}, \quad \text{with} \quad \beta = \frac{1}{RT} \\ &= \frac{1}{Q} \sum_{s|F \in s} e^{-\beta E(s)}, \quad \text{with} \quad Q = \sum_s e^{-\beta E(s)} \end{aligned}$$

Efficient dynamic programming algorithms are available to compute

- Minimum free energy $MFE = \min_s E(s)$
- Partition function Q
- Probabilities p_{kl} for base pairs (k, l) , i.e. $p_{kl} = p(F)$ with $F = (k, l)$
- Suboptimal and locally stable structures, ligand binding, ...

Implementations for different input data types^[2]

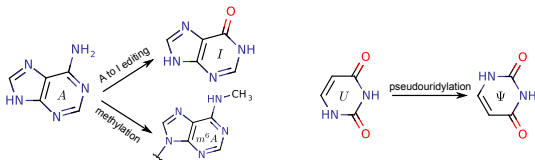
- single sequences
- (multiple) interacting sequences
- sequence alignments (consensus structure)

Implementations are restricted to unmodified bases!

[2] R. Lorenz et al. "ViennaRNA Package 2.0". In: *Algorithms for Molecular Biology* 6 (2011), pp. 1–14. DOI: 10.1186/1748-7188-6-26

Modified Bases in RNA

Post-transcriptional RNA modifications (epitranscriptome):

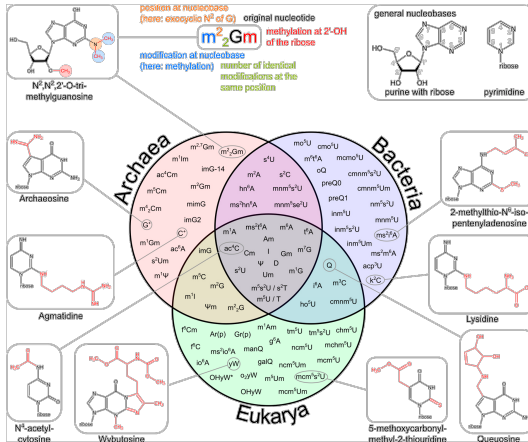


- Modomics Database^[3]: > 170 distinct modified bases
- Well known modifications: *I*, Ψ , m^6A , m^1A , m^5C , ...
- Function and purpose of modifications still largely unknown
- Structural effects of base modifications:
 - ① correct folding into functional structures (tRNA, rRNA, etc.)
 - ② regulation of protein binding sites (mRNAs, lncRNAs)
 - ③ regulation of RNA-RNA binding sites (siRNA, miRNA)
 - ④ Modifications may change pairing partner preference
 - ⑤ Modifications may (de-)stabilize loop formation

[3] A. Cappannini et al. "MODOMICS: a database of RNA modifications and related information. 2023 update". In: *Nucleic Acids Research* 52.D1 (2024), pp. D239–D244. DOI: 10.1093/nar/gkad1083

Modifications in tRNA^[4]

- 93 known post-transcriptional modifications



- Modifications can be subtle from the RNA structure perspective
- Some are essential to induce structural domain rearrangements

[4] C. Lorenz et al. "tRNA modifications: impact on structure and thermal adaptation". In: *Biomolecules* 7.2 (2017), p. 35. DOI: 10.3390/biom7020035

RNA Secondary Structure Prediction and Modified Bases

How to incorporate modified bases in predictions?

Obstacles:

- 2D structure effects known only for a minority of modifications
- 3D effects may yield non-Nearest-Neighbor dependencies
- Combinatorial explosion for energy parameters and pairing rules

Possible Solutions:

- Re-implement algorithms with
 - ① Enhanced nucleotide alphabet
 - ② Additional pairing rules
 - ③ More energy parameters
- Use constraints^[5]
 - ① Hard constraints: pairing rules
 - ② Soft constraints: supplement (partial) energy parameters

[5] R. Lorenz et al. "RNA folding with hard and soft constraints". In: *Algorithms for Molecular Biology* 11.1 (2016), pp. 1–13. DOI: 10.1186/s13015-016-0070-z

RNA Secondary Structure Prediction and Modified Bases

How to incorporate modified bases in predictions?

Obstacles:

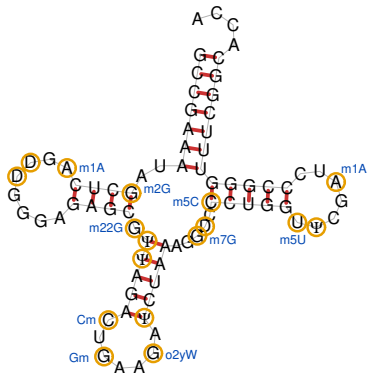
- 2D structure effects known only for a minority of modifications
- 3D effects may yield non-Nearest-Neighbor dependencies
- Combinatorial explosion for energy parameters and pairing rules

Possible Solutions:

- Re-implement algorithms with
 - ① Enhanced nucleotide alphabet
 - ② Additional pairing rules
 - ③ More energy parameters
- **Use constraints**^[5]
 - ① Hard constraints: pairing rules
 - ② Soft constraints: supplement (partial) energy parameters

[5] R. Lorenz et al. "RNA folding with hard and soft constraints". In: *Algorithms for Molecular Biology* 11.1 (2016), pp. 1–13. DOI: 10.1186/s13015-016-0070-z

Example: human tRNA Phe (tdbR00000103)^[6]

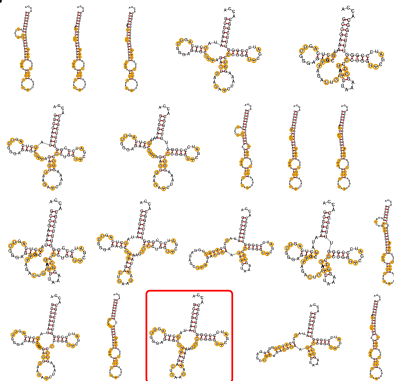
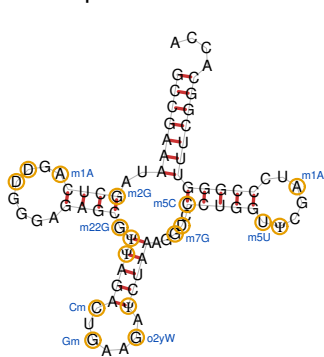


- 17 out of 76 nucleotides are modified
- Predicted ground state (MFE) is **not** cloverleaf structure

[6] F. Jühling et al. "tRNAdb 2009: compilation of tRNA sequences and tRNA genes". In: *Nucleic Acids Research* 37.suppl.1 (2009), pp. D159–D162

Example: human tRNA Phe (tdbR00000103)

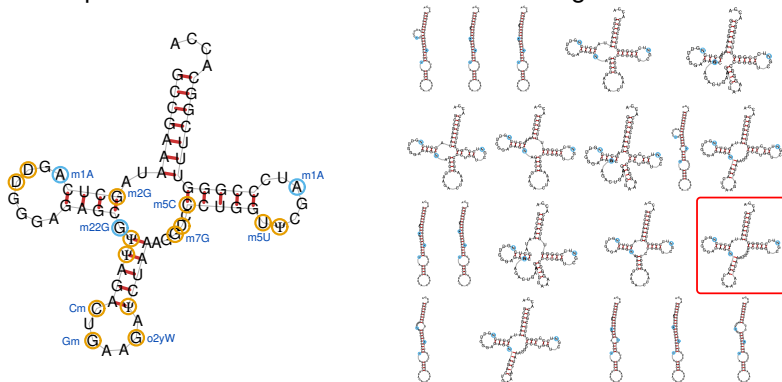
1. Suboptimals for *unmodified* sequence



- Correct cloverleaf structure at position 18
- $\bar{d}_{BP} = \frac{1}{|S^3|} \sum_{s \in S^3} d_{BP}(s_{ref}, s) = 21.55 \text{ bp}$, $E_{ref} - E_{MFE} = 1.7 \text{ kcal/mol}$
- Most predictions are not cloverleaf shaped

Example: human tRNA Phe (tdbR00000103)

3. Suboptimals with *masked bases* for RT blocking modifications^[7]



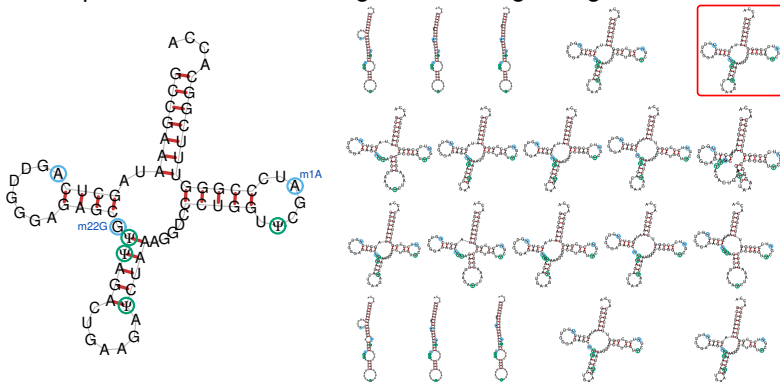
- Correct cloverleaf structure at position 15
- $\bar{d}_{BP}^3 = 19.7$ bp, $E_{ref} - E_{MFE} = 1.6$ kcal/mol
- Still, most predicted structure are not tRNA-like!

Ψ-A stacking energies are available in literature!

[7] Y. Motorin et al. "Identification of modified residues in RNAs by reverse transcription-based methods". In: *Methods in Enzymology* 425 (2007), pp. 21–53. DOI: 10. 1016/S0076-6879(07)25002-5

Example: human tRNA Phe (tdbR00000103)

4. Suboptimals with RT masking and stacking energies for Ψ -A^[8]

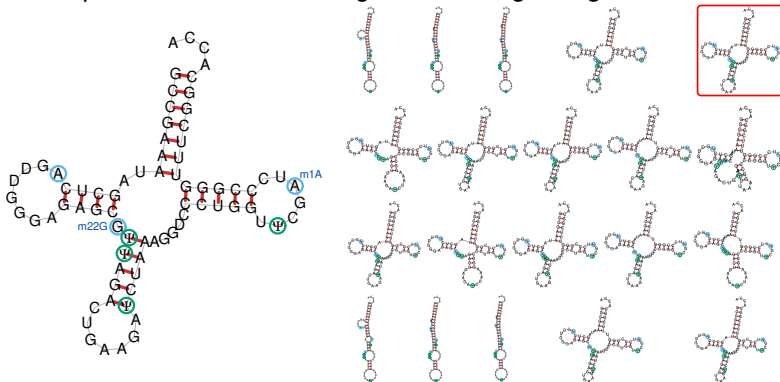


- Correct cloverleaf structure at position 5
- $\bar{d}_{BP}^3 = 15.5$ bp, $E_{ref} - E_{MFE} = 0.83$ kcal/mol
- Can we do better?

[8] G. A. Hudson et al. "Thermodynamic contribution and nearest-neighbor parameters of pseudouridine-adenosine base pairs in oligoribonucleotides". In: *RNA* 19.11 (2013), pp. 1474–1482. DOI: 10.1261/rna.039610.113

Example: human tRNA Phe (tdbR00000103)

4. Suboptimals with RT masking and stacking energies for Ψ -A^[8]



- Correct cloverleaf structure at position 5
- $\bar{d}_{BP}^3 = 15.5$ bp, $E_{ref} - E_{MFE} = 0.83$ kcal/mol
- Can we do better?

Dihydrouridine (D) adds flexibility to the D-loop!

[8] G. A. Hudson et al. "Thermodynamic contribution and nearest-neighbor parameters of pseudouridine-adenosine base pairs in oligoribonucleotides". In: *RNA* 19.11 (2013), pp. 1474–1482. DOI: 10.1261/rna.039610.113

Example: human tRNA Phe (tdbR00000103)

No (stacking) energies with dihydrouridine (D) are available

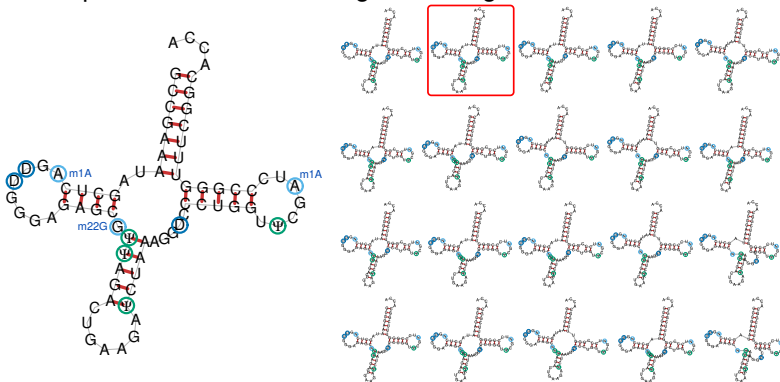
Structural effects of dihydrouridine^[9]

- C3'-endo sugar conformation is destabilized in favor of C2'-endo
- more flexibility
- promotes destacking
- destabilization is somewhere in the range of 1.5 kcal/mol

[9] J. J. Dalluge et al. "Conformational flexibility in RNA: the role of dihydrouridine". In: *Nucleic Acids Research* 24.6 (1996), pp. 1073–1079. DOI: 10.1093/nar/24.6.1073

Example: human tRNA Phe (tdbR00000103)

5. Suboptimals with RT masking and energies for Ψ -A and D^[9]

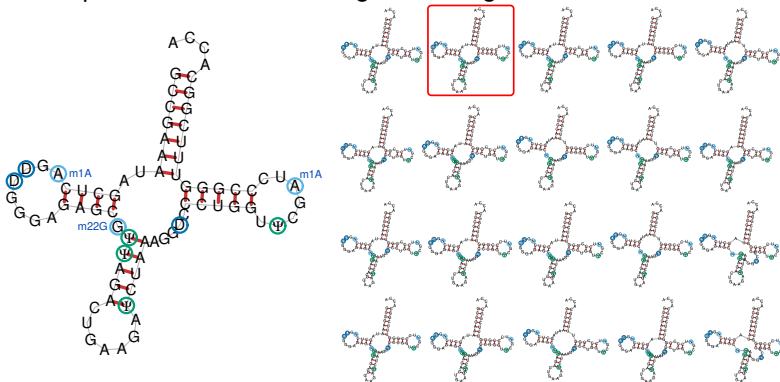


- Correct cloverleaf structure at position 2
- $\bar{d}_{BP}^3 = 6.92$ bp, $E_{ref} - E_{MFE} = 0.2$ kcal/mol

[9] J. J. Dalluge et al. "Conformational flexibility in RNA: the role of dihydrouridine". In: *Nucleic Acids Research* 24.6 (1996), pp. 1073–1079. DOI: 10. 1093/nar/24. 6. 1073

Example: human tRNA Phe (tdbR00000103)

5. Suboptimals with RT masking and energies for Ψ -A and D^[9]



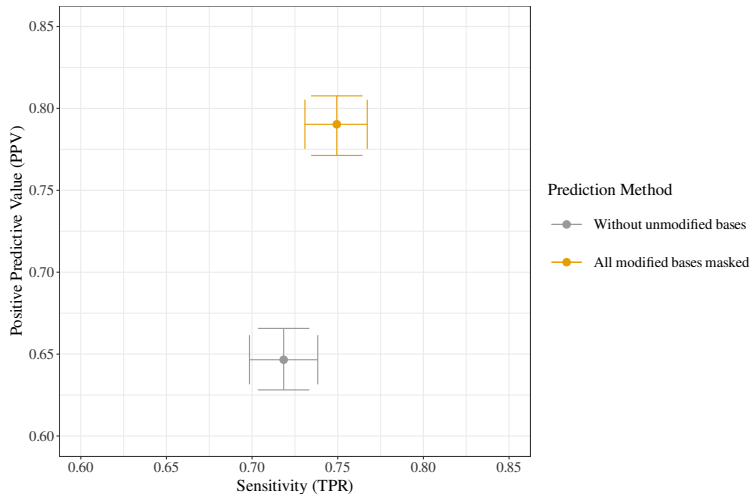
- Correct cloverleaf structure at position 2
- $\bar{d}_{BP}^3 = 6.92$ bp, $E_{ref} - E_{MFE} = 0.2$ kcal/mol

How do these constraints affect prediction for other tRNAs?

[9] J. J. Dalluge et al. "Conformational flexibility in RNA: the role of dihydrouridine". In: *Nucleic Acids Research* 24.6 (1996), pp. 1073–1079. DOI: 10. 1093/nar/24. 6. 1073

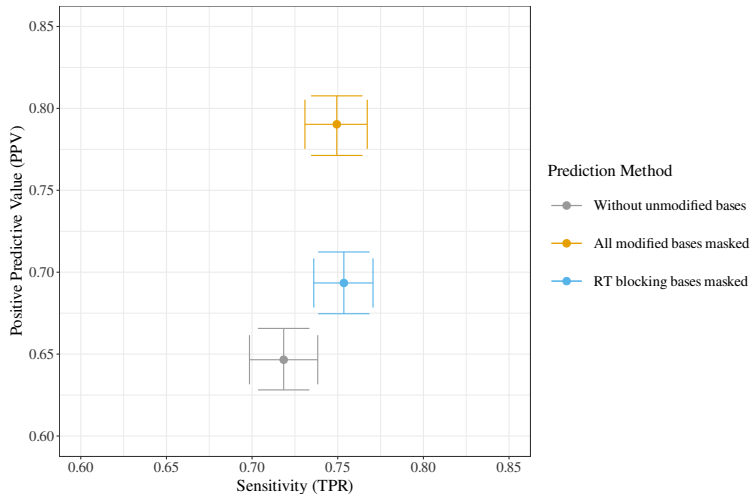
Overall prediction improvement

Performance on tRNADB data set (623 sequences)



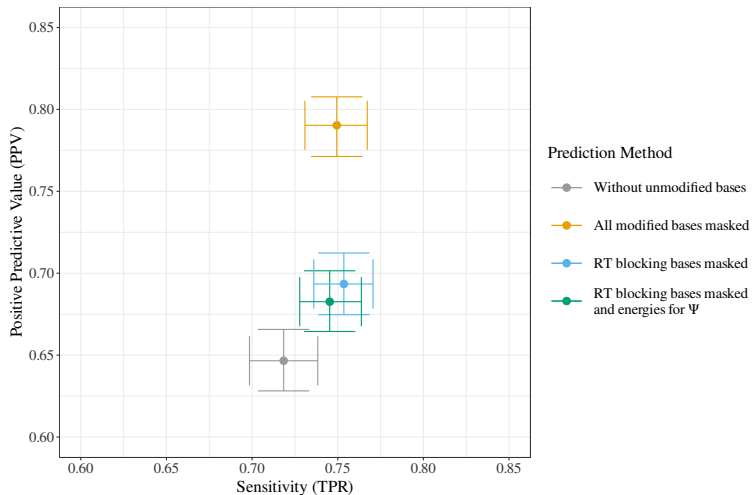
Overall prediction improvement

Performance on tRNADB data set (623 sequences)



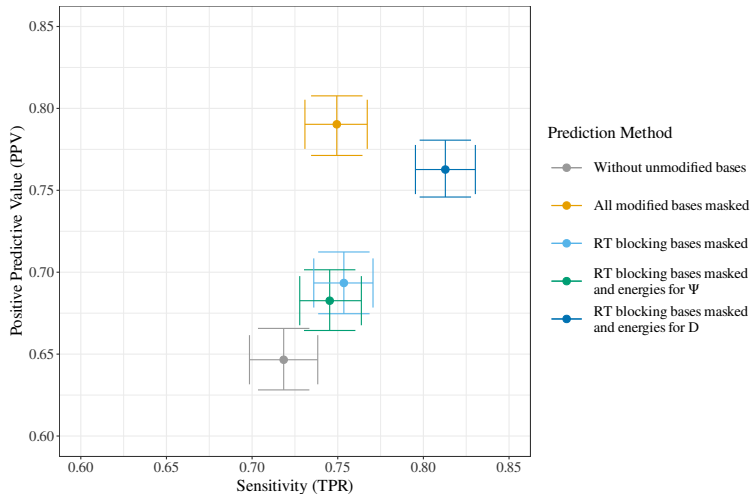
Overall prediction improvement

Performance on tRNADB data set (623 sequences)



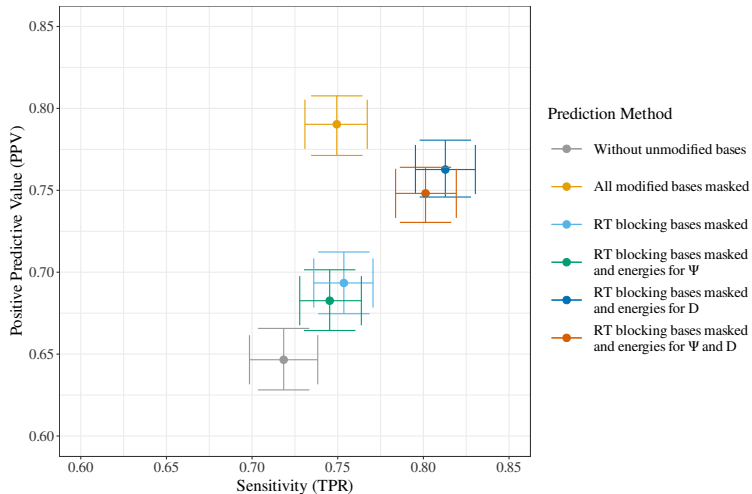
Overall prediction improvement

Performance on tRNAdb data set (623 sequences)



Overall prediction improvement

Performance on tRNADB data set (623 sequences)



Secondary Structures with Modified Bases

Implemented into the ViennaRNA Package^[10]:

- Includes stacking parameters:

One-letter-code	Modified base
7	7-deaza-adenosine (7DA)
I	Inosine
6	N6-methyladenosine (m^6A)
P	Pseudouridine
9	Purine (a.k.a. nebularine)
D	Dihydrouridine (Rosetta-RECESS)

- Extension for arbitrarily many additional modifications
- Available from command line and API
- Easy extension with more data using JSON parameter files (stacking, terminal, mismatch, dangles, fallback, pairing partners)

[10] Y. Varenik et al. "Modified RNAs and predictions with the ViennaRNA Package". In: *Bioinformatics* 39.11 (2023), btad696. DOI: 10.1093/bioinformatics/btad696

Secondary Structures with Modified Bases

Implemented into the ViennaRNA Package^[10]:

- Includes stacking parameters:

One-letter-code	Modified base
7	7-deaza-adenosine (7DA)
I	Inosine
6	N6-methyladenosine (m^6A)
P	Pseudouridine
9	Purine (a.k.a. nebularine)
D	Dihydrouridine (Rosetta-RECESS)

- Extension for arbitrarily many additional modifications
- Available from command line and API
- Easy extension with more data using JSON parameter files (stacking, terminal, mismatch, dangles, fallback, pairing partners)

Sometimes including modified base yields worse predictions!

[10] Y. Varenik et al. "Modified RNAs and predictions with the ViennaRNA Package". In: *Bioinformatics* 39.11 (2023), btad696. DOI: 10.1093/bioinformatics/btad696

Outlook

**Typically parameters are obtained from UV-melting experiments
(expensive)**

Outlook

Typically parameters are obtained from UV-melting experiments (expensive)

Parameters:

- Add (more) in-silico predicted parameters^{[11][12]}
- Add parameters obtained from fluorescence melting/titration^[13]
- Add modifications known to stay unpaired (literature search, experts knowledge)

Implementation:

- Comparative structure prediction
- Multiple interacting RNAs
- Pairs between different modified bases

[11] F.-C. Chou et al. "Blind tests of RNA nearest-neighbor energy prediction". In: *Proceedings of the National Academy of Sciences* 113.30 (2016), pp. 8430–8435. DOI: 10.1073/pnas.152333511

[12] M. C. Hopfinger et al. "Predictions and analyses of RNA nearest neighbor parameters for modified nucleotides". In: *Nucleic Acids Research* 48.16 (2020), pp. 8901–8913. DOI: 10.1093/nar/gkaa654

[13] J. Wang et al. "Assessment for melting temperature measurement of nucleic acid by HRM". In: *Journal of analytical methods in chemistry* 2016.1 (2016), p. 5318935. DOI: 10.1155/2016/5318935

Acknowledgements

TBI Vienna

- Yuliia Varenyk
- Thomas Spicher
- Ivo L. Hofacker

MedUni Vienna

- Michael Jantsch
- Hamid Reza Mansouri Khosravi

Uni Leipzig

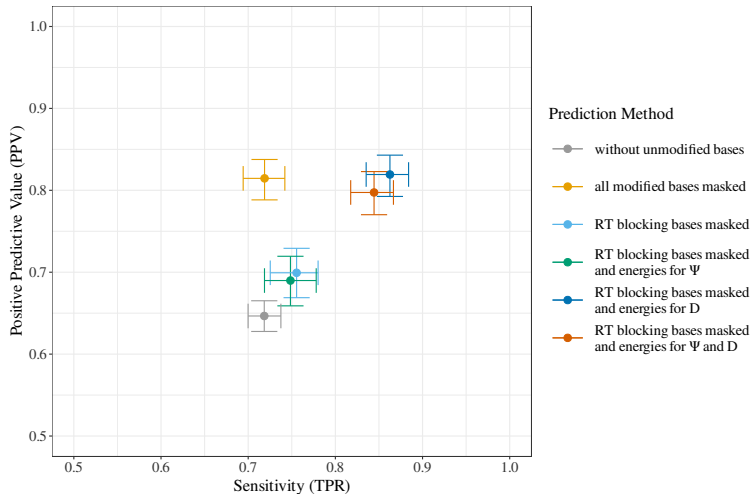
- Peter F. Stadler + group

Thank You for your attention!



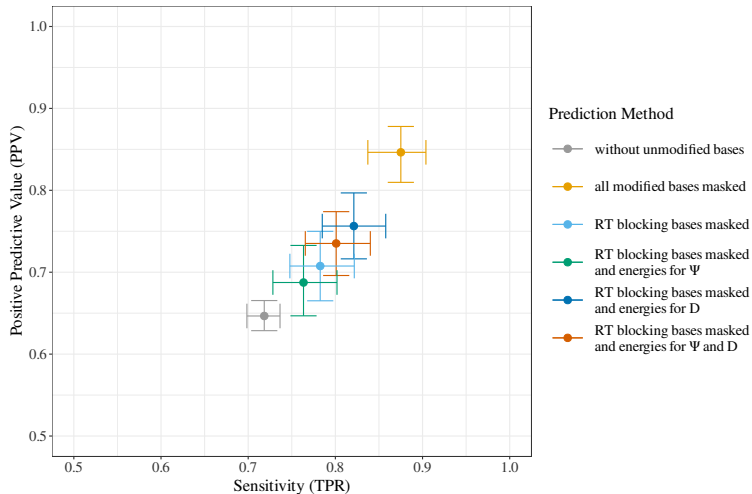
tRNA Secondary Structure Prediction

Performance on tRNADB data set (eucaryotes, 242 sequences)



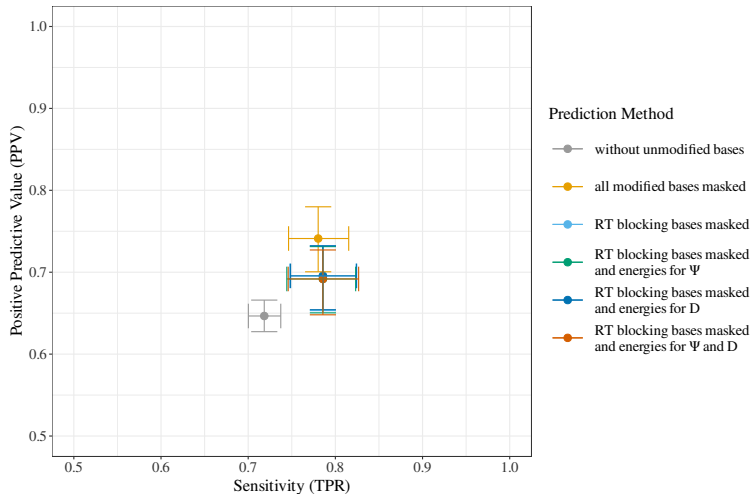
tRNA Secondary Structure Prediction

Performance on tRNADB data set (bacteria, 139 sequences)



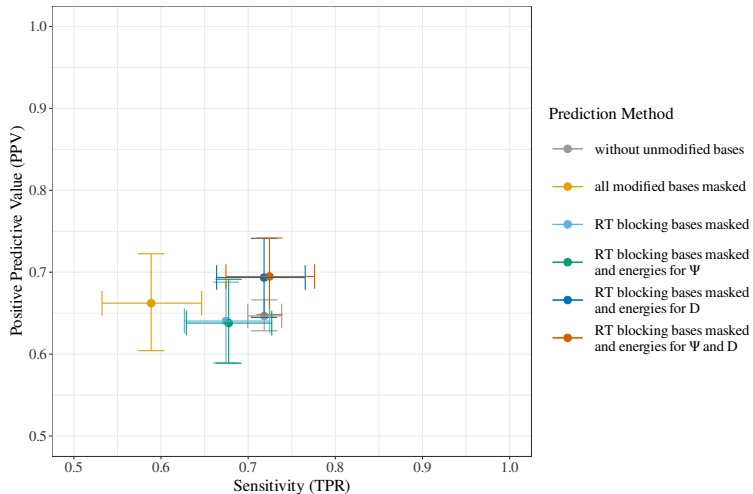
tRNA Secondary Structure Prediction

Performance on tRNADB data set (archaea, 76 sequences)



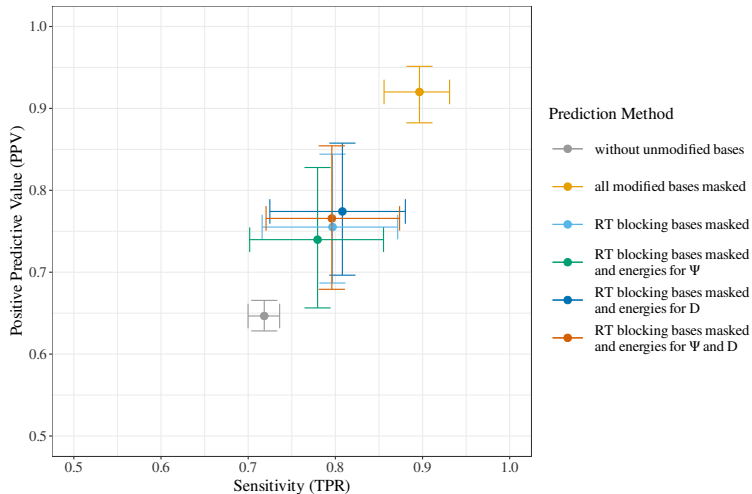
tRNA Secondary Structure Prediction

Performance on tRNADB data set (eucaryotes_mito, 111 sequences)



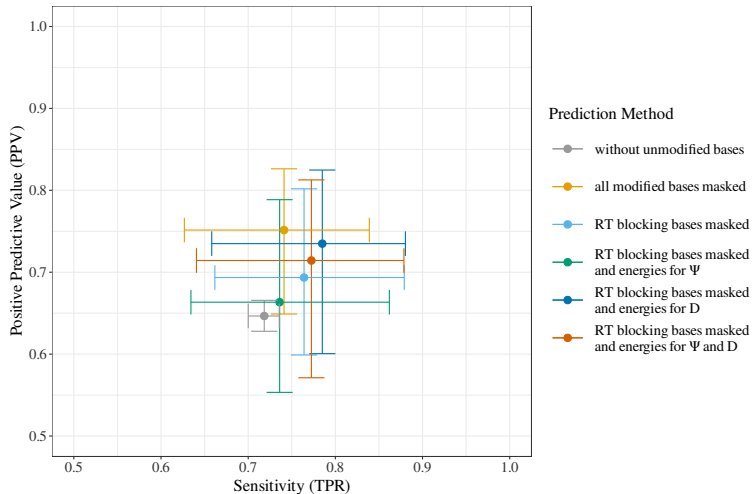
tRNA Secondary Structure Prediction

Performance on tRNADB data set (eucaryotes_plastids, 38 sequences)



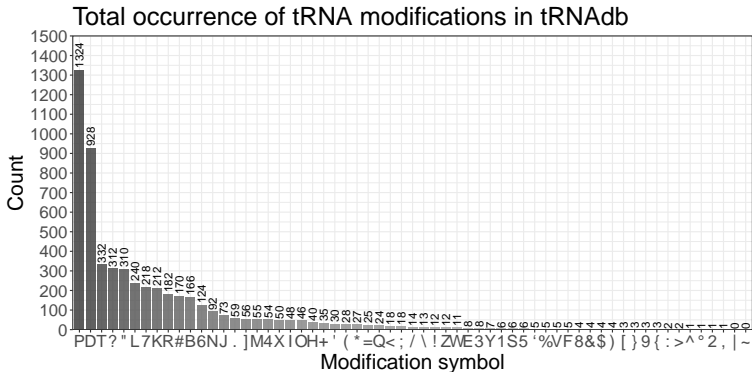
tRNA Secondary Structure Prediction

Performance on tRNAdb data set (virus, 17 sequences)



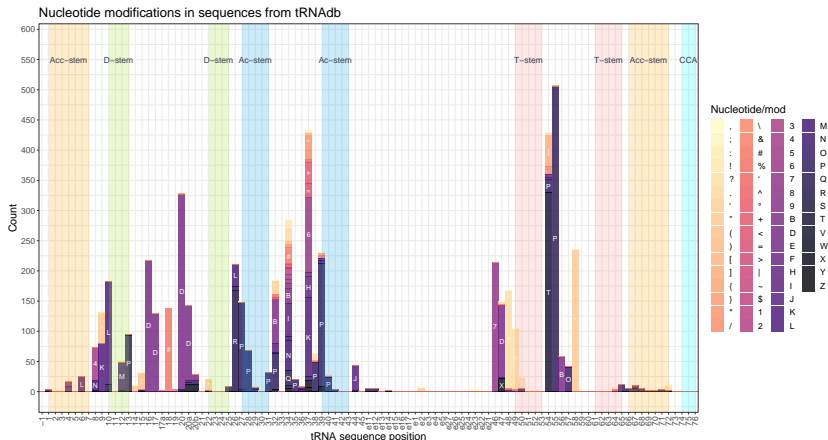
Modifications in tRNA

- Frequency of modifications in tRNAdb



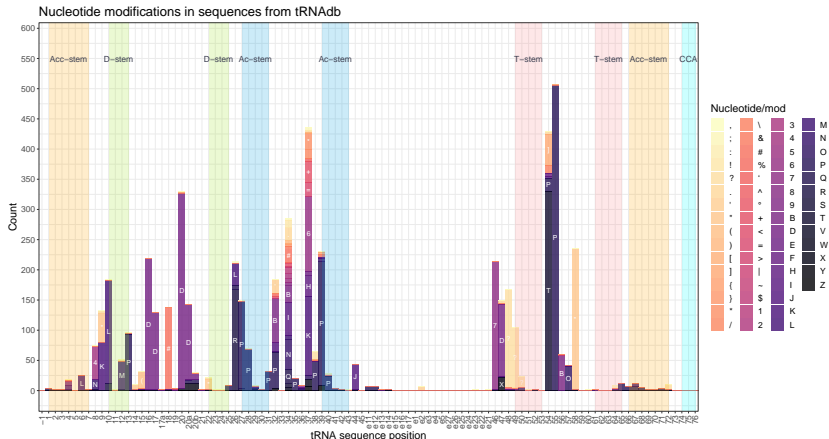
Modifications in tRNA

- Frequency of modifications in tRNAdb
- Which modifications can be found where?



Modifications in tRNA

- Frequency of modifications in tRNAdb
- Which modifications can be found where?
- Which modifications might induce structural rearrangements?



Nearest Neighbor Parameters for Modified Bases

Base pair stacking and helix ends

- I-U (Wright et al. 2007, partial, UV-melting)
- I-U (Jolley et al. 2015, partial, in-silico)
- I-C, D-U, iG-iC (Chou et al. 2016, partial, in-silico)
- I-C, iG-iC, D-U, I-U (Hopfinger et al. 2020, partial, in-silico)
- I-C (Wright et al. 2018, partial, UV-melting)
- Ψ -A (Hudson et al. 2013, partial, UV-melting)
- Ψ -A (Deb et al. 2019, partial, in-silico, NMR)
- Ψ -{A, U, C, G} (Kierzek et al. 2014, partial, UV-melting, CD, NMR)
- m^6A -U (Roost et al. 2015, partial, UV-melting, NMR)
- modified U -{A, U, C, G}, I -{A, C, U} (Vendeix et al. 2009, partial, in-silico)
- P-U (Jolley et al. 2017, partial, UV-melting)
- 7DA-U (Richardson et al. 2016, partial, UV-melting)

Other data on structural effects

- Dihydrouridine (Dalluge et al. 1996, Dyubankova et al. 2015, NMR)
- 2-thiouridine (s^2U), Ψ , dihydrouridine (Sipa et al. 2007, partial, CD)
- tRNA modifications (Lorenz C. et al. 2017, meta)
- m^6A , m^1A , m^5C , Ψ , I (Harcourt et al. 2017, meta)