Modified Nucleotides in RNA structure prediction

Ronny Lorenz, Thomas Spicher, Ivo L. Hofacker ronny@tbi.univie.ac.at

University of Vienna, Theoretical Biochemistry Group (TBI)

Mondsee, Austria, September 8th 2020







RNA Secondary Structures RNAs fold hierarchical (A) \rightarrow (B) \rightarrow (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
- Starting point for 3D structure modeling
- Good abstraction to perform folding kinetics simulation

RNA Secondary Structures - Nearest Neighbor Energy Model



- 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_L¹

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

¹Turner et al., "NNDB: The nearest neighbor parameter database for predicting stability of nucleic acid secondary structure.", 2009, NAR 38, D280–D282

RNA Secondary Structures - Statistical Thermodynamics

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

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, ...

The ViennaRNA $Package^2$ provides fast implementations for

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

²Lorenz et al., "ViennaRNA Package 2.0.", 2011, ALMOB 6.1, 26

RNA Secondary Structures - Statistical Thermodynamics

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

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, ...

The ViennaRNA $Package^2$ provides fast implementations for

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

Implementations are restricted to unmodified bases!

²Lorenz et al., "ViennaRNA Package 2.0.", 2011, ALMOB 6.1, 26

RNA Secondary Structure - Modified Bases

Many RNAs require modified bases to fulfill their function!

- Modifications may change pairing partner preference
- Modifications may (de-)stabilize loop formation
- Modomics Database³ lists 172 different modified bases
- Only a few energy parameters involving modified bases available, e.g. stacking with ΨA^4 and $I U^5$ pairs

How to incorporate modified bases in predictions?

- Use constraints
 - Mask modified bases, e.g. to stay unpaired using letter "N"
 - Write wrapper (Python/Perl/C/C++) using the constraints framework of the ViennaRNA Package⁶ to supplement energy parameters
- Re-implement algorithms with
 - enhanced nucleotide alphabet
 - additional pairing rules
 - 8 more energy parameters

³Boccaletto et al., "MODOMICS: a database of RNA modification pathways. 2017 update.", 2018, NAR 46.D1, D303–D307

⁴Hudson et al., "Thermodynamic contribution and nearest-neighbor parameters of pseudouridine-adenosine base pairs in oligoribonucleotides,", 2013, RNA 19.11, 1474–1482

⁵Wright et al., "Nearest Neighbor Parameters for Inosine-Uridine Pairs in RNA Duplexes.", 2007, Biochemistry 46, 4625–4634 ⁶Lorenz et al., "RNA folding with hard and soft constraints.", 2016, ALMOB 11, 8

RNA Secondary Structure - Modified Bases

Many RNAs require modified bases to fulfill their function!

- Modifications may change pairing partner preference
- Modifications may (de-)stabilize loop formation
- Modomics Database³ lists 172 different modified bases
- Only a few energy parameters involving modified bases available, e.g. stacking with ΨA^4 and $I U^5$ pairs

How to incorporate modified bases in predictions?

Use constraints

- Mask modified bases, e.g. to stay unpaired using letter "N"
- Write wrapper (Python/Perl/C/C++) using the constraints framework of the ViennaRNA Package⁶ to supplement energy parameters
- Re-implement algorithms with
 - enhanced nucleotide alphabet
 - 2 additional pairing rules
 - 8 more energy parameters

³Boccaletto et al., "MODOMICS: a database of RNA modification pathways. 2017 update.", 2018, NAR 46.D1, D303–D307

⁴Hudson et al., "Thermodynamic contribution and nearest-neighbor parameters of pseudouridine-adenosine base pairs in oligoribonucleotides,", 2013, RNA 19.11, 1474–1482

⁵Wright et al., "Nearest Neighbor Parameters for Inosine-Uridine Pairs in RNA Duplexes.", 2007, Biochemistry 46, 4625–4634 ⁶Lorenz et al., "RNA folding with hard and soft constraints.", 2016, ALMOB 11, 8

Example: human tRNA Phe (tdbR00000103⁷)



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

⁷Jühling et al., "tRNAdb 2009: compilation of tRNA sequences and tRNA genes.", 2009, NAR 37, D159–D162

Example: human tRNA Phe (tdbR00000103⁷)



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

What about suboptimal structures $S^{\delta} = \{s \in \Omega \mid E(s) \le E_{\mathsf{MFE}} + \delta\}$?

⁷Jühling et al., "tRNAdb 2009: compilation of tRNA sequences and tRNA genes.", 2009, NAR 37, D159–D162

1. Suboptimals for *unmodified sequence*



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

1. Suboptimals for *unmodified sequence*



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

How about masking modified bases?

2. Suboptimals for sequence with all modified bases masked



- Modified bases are excluded from pairing
- · Correct cloverleaf structure not among the solution set
- $\bar{d}_{BP}^3 = 16.7$ bp, anticodon arm is missing

2. Suboptimals for sequence with all modified bases masked





- Modified bases are excluded from pairing
- · Correct cloverleaf structure not among the solution set
- $\bar{d}_{BP}^3 = 16.7$ bp, anticodon arm is missing

Some modified bases are required to pair!

3. Suboptimals with masked bases for RT blocking modifications⁸



- 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!

⁸Motorin et al., "Identification of modified residues in RNAs by reverse transcription-based methods.", 2007, Methods in Enzymology 425, 21–53

3. Suboptimals with masked bases for RT blocking modifications⁸



- 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!

⁸Motorin et al., "Identification of modified residues in RNAs by reverse transcription-based methods.", 2007, Methods in Enzymology 425, 21–53

Example: human tRNA Phe (tdbR00000103) Python overhead to add energies for stacks with Ψ -A⁹:

Energy parameters from Hudson et al., 2013

```
stacking_psi_A = [
```

	# N	Α	С	G	U	Р	
	[Ο,	Ο,	Ο,	Ο,	Ο,	0],	# N
	[Ο,	Ο,	Ο,	Ο,	-77,	0],	# A
	[Ο,	Ο,	Ο,	-14,	Ο,	0],	# C
	[Ο,	Ο,	-9,	Ο,	Ο,	0],	# G
	[Ο,	-181,	Ο,	Ο,	Ο,	0],	# U
	[Ο,	Ο,	Ο,	Ο,	Ο,	0]	# F
1							

stacking_A_psi = [

1

# N	Α	С	G	U	Р	
[Ο,	Ο,	Ο,	Ο,	Ο,	0],	# 1
[Ο,	Ο,	Ο,	Ο,	-69,	0],	# 1
[Ο,	Ο,	Ο,	-105,	Ο,	0],	# (
[Ο,	Ο,	-69,	Ο,	Ο,	0],	# (
[Ο,	-170,	Ο,	Ο,	Ο,	0],	# l
[Ο,	Ο,	Ο,	Ο,	Ο,	0]	# 1

```
# Energy correction callback passed to ViennaRMA library
def correction(i, j, k, l, step, encoding):
    if (step == RNA.DECOMP_PAIR_IL): # an internal loop is evaluated
        if i + 1 == k and j - 1 == 1: # (i,j) and (k, l) are stacking
        if (encoding[i] == 5 and encoding[j] == 1):
            return stacking_psi_A[encoding[i] == 5):
            return stacking_psi_A[encoding[i] == 5):
            return stacking_si_A[encoding[i] == 5):
            return stacking_Apsi_encoding[i] == 5):
            return stacking_Apsi[encoding[i] == 5):
            return stacking_Apsi[encoding[i] == 1):
            return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
         return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i]] == 1):
        return stacking_Apsi[encoding[i] == 1):
        return stacking_Apsi[encoding[i] == 1):
        return stacking_Apsi[encoding[i] == 1):
         return stacking_Apsi[encoding[i] == 1):
         return stacking_Apsi[enc
```

return 0

⁹Hudson et al., "Thermodynamic contribution and nearest-neighbor parameters of pseudouridine-adenosine base pairs in oligoribonucleotides", 2013, RNA 19.11, 1474–1482

4. Suboptimals with RT masking and stacking energies for Ψ -A¹⁰



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

¹⁰Hudson et al., "Thermodynamic contribution and nearest-neighbor parameters of pseudouridine-adenosine base pairs in oligoribonucleotides,", 2013, RNA 19.11, 1474–1482

4. Suboptimals with RT masking and stacking energies for Ψ -A¹⁰



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

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

¹⁰Hudson et al., "Thermodynamic contribution and nearest-neighbor parameters of pseudouridine-adenosine base pairs in oligoribonucleotides.", 2013, RNA 19.11, 1474–1482

Example: human tRNA Phe (tdbR00000103) No (stacking) energies with dihydrouridine (D) are available

Structural effects of dihydrouridine¹¹

- C3'-endo sugar conformation is destabilized in favor of C2'-endo
- more flexibility

return e

- promotes destacking
- destabilization of 1.5 kcal/mol (mono), up to 5.3 kcal/mol (oligo)

Required changes for correction function:

```
def correction(i, j, k, l, step, encoding):
    e = 0 # default correction
    if (step == RNA.DECOMP_PAIR_IL): # an internal loop is evaluated
        if \mathbf{i} + 1 = \mathbf{k} and \mathbf{i} - 1 = \mathbf{l}; # (i, i) and (k, l) are stacking
            if (encoding[i] == 5 and encoding[i] == 1);
                e += stacking_psi_A[encoding[k]][encoding[1]]
            elif (encoding[k] == 1 and encoding[1] == 5):
                e += stacking_psi_A[encoding[j]][encoding[i]]
            elif (encoding[i] == 1 and encoding[j] == 5):
                 e += stacking_A_psi[encoding[k]][encoding[1]]
            elif (encoding[k] == 5 and encoding[1] == 1):
                e += stacking_A_psi[encoding[j]][encoding[i]]
            if encoding[i] == 6 or \
                encoding[j] == 6 or \
                encoding[k] == 6 or \
                 encoding[1] == 6:
                e += 150 # D destabilizes stacking by at least 1.5 kcal/mol
```

¹¹ Dalluge et al., "Conformational flexibility in RNA: the role of dihydrouridine.", 1996, NAR 24.6, 1073–1079

5. Suboptimals with RT masking and energies for Ψ -A and D¹²



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

¹² Dalluge et al., "Conformational flexibility in RNA: the role of dihydrouridine.", 1996, NAR 24.6, 1073-1079

5. Suboptimals with RT masking and energies for Ψ -A and D¹²



- 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?

¹² Dalluge et al., "Conformational flexibility in RNA: the role of dihydrouridine.", 1996, NAR 24.6, 1073–1079











Performance on tRNAdb data set (623 sequences)

Conclusion Modified bases may heavily influence structure space

The good news

- Some data on structural effects of modifications available
- Better prediction performance through constraints
- Constraints available for many prediction algorithms

The not so good news

- Additional parameters do not necessarily increase performance
- Constraints become complex for more modifications and contexts
- Unrealistic to include full parameters with many modified bases
- No unique base annotation (tRNAdb¹³, RNAmod¹⁴, $\underline{MODOMICS}^{15}$)

Outlook

- Gather more data on structural effects of modified bases
- Rule and energy parameter set for pairs with modified bases
- Integration of modified bases in ViennaRNA Package

¹³ Jühling et al., "tRNAdb 2009: compilation of tRNA sequences and tRNA genes.", 2009, NAR 37, D159–D162
¹⁴ Liu et al., "RNAmod: an integrated system for the annotation of mRNA modifications", 2019, NAR 47, W548-W555
¹⁵ Boccaletto et al., "MODOMICS: a database of RNA modification pathways. 2017 update", 2018, NAR 46, D303-D307

Acknowledgements

- Ivo L. Hofacker
- Thomas Spicher
- Peter F. Stadler
- Yuliia Varenyk
- The remaining TBI group

Thank You for your attention!















Performance on tRNAdb data set (eucaryotes_plastids, 38 sequences)

