

Structural Bioinformatics Final Exam

RNA part report

Antonio Ortega Jiménez

January 27, 2017

Exercise 1: RNAfold prediction on seq1

```
AUCAGUUCUAGCAGGAGCUGUACUCAGAGACUCGGGAAAUUUUCCGGAAUUUUACCCGGGUUUUUACGU
..(((((((.....)))))).....(((((((((((.....(((.....)))..)))))))))).....
```

a) The obtained Minimum Free Energy (MFE) secondary structure consists of a multi-loop with two closed components. The left arm contains a hairpin, whereas the right arm contains another hairpin with a couple of internal loops. The centroid version discards these internal loops and instead shows a simple hairpin with a huge loop. This is due to the low probability exhibited by the base pairs in the terminal region of the hairpin [?].

b) RNAfold predicts 18 base pairs with a probability higher than 0.8.

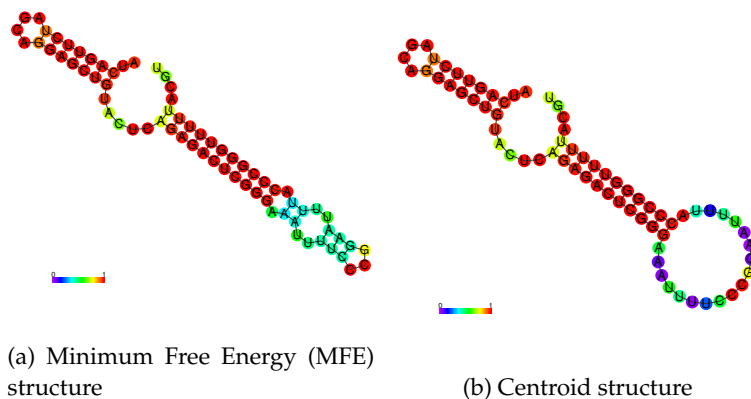


Figure 1: Structures predicted by RNAFold for the query sequence. The color code shows base pair probabilities.

Exercise 2: RNAfold prediction on seq2

```
AUCGGUUCACAGCAGGAACUGUACUCGGGGGCUCGGGAAACCCUCCGGGGUUUUACCCGGGUUUUUACGU
..(((((((.....)))))).....(((((((((((.....(((.....)))..)))))))))).....
```

a) *A priori*, this sequence has a similar structure. Base pairs are distributed differently in the right arm, and this has triggered an increase in the hairpin stability. This can be easily noticed:

- The centroid structure now includes the hairpin, because there are more structures in the ensemble that now feature the same hairpin.

- Instead of a highly stable duplex and a highly unstable loop, stability is uniformly distributed across the whole hairpin (see second bulge in right violin from figure 4). This enables the formation of the most interior bonds in the duplex, right before the hairpin's loop, which therefore grows smaller.
- Thermodynamical ensemble prediction parameters show that the structural diversity is reduced and the frequency of the MFE structure increases.

b) RNAfold predicts 20 base pairs with a probability higher than 0.8.

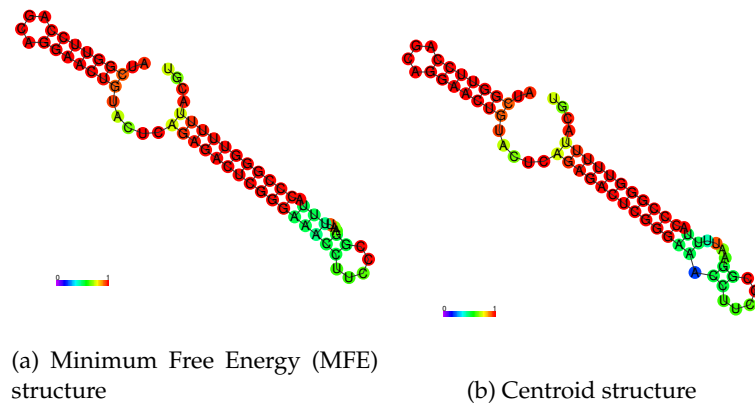


Figure 2: Structure predicted by RNAfold for the query sequence. The color code shows base pair probabilities.

Exercise 3: Hamming distance

Hamming distance between structures is 22, and just 11 at the sequence level.

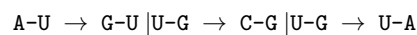
Exercise 4: Base pair distance

Base pair distance between structures is 30.

Exercise 5: Discussion

Both sequences have a high identity, with a Hamming distance of just 11 at the sequence level. Nevertheless, these differences are not uniformly distributed across the sequence, and thus, they exert different impacts on each arm.

- In spite of the 3 mutations suffered by the left arm, on positions 3 10 and 18, the structure is not changed. This is due to the consistent nature of these mutations, that preserve the affected base pairs:



The dotplots (figure 3) support the idea that the left arm has not changed: no significant differences are observed in the left arm area (top left corner), and no alternative structures are observed within the probability space.

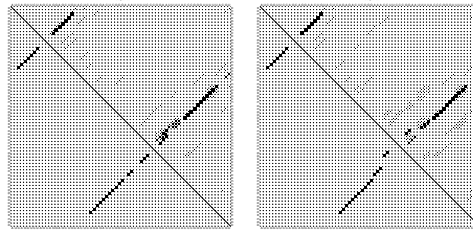


Figure 3: Base pair probability dot plot for seq1 (left) and seq2 (right)

- The remaining 8 mutations are located on the right arm, and though the structural motif is preserved, its internal organization has changed totally.

On the one hand, base pairs cluster into 2 blocks of the same length (7 base pairs each, 14 in total), instead of 3 blocks of length 11, 2 and 3 (16 in total). This organization leads to a unique internal loop in the duplex center. Since internal loops and bulges destabilize duplexes, as stated in [] it follows that a reorganization that reduces the number of internal loops is bound to make the structure more stable.

```
..(((((((.....)))))).....(((((((((((.....(((.....)))..)))))))))).....
..(((((((.....))))))..(((((((.....(((((((.....)))))).....)))))).....
```

left arm

right arm

On the other hand, this reorganization leads to a redistribution of base pairing probability toward inner regions of the duplex, that become more stable. In other words, in exchange for the loss of 2 base pairs, the remaining ones become more stable and trigger the stabilization of the hairpin's duplex. As a result, the whole hairpin in the right arm is more probable. This phenomenon can be easily observed in figure 4.

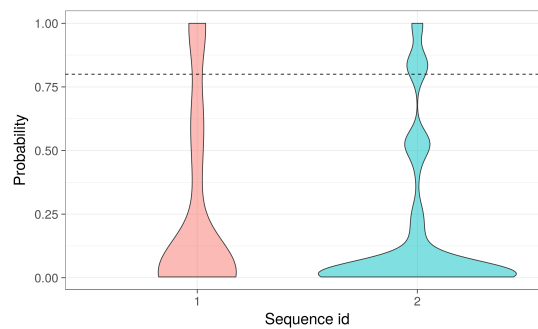


Figure 4: Violin plot showing the base pair probabilities distribution. In sequence 1, we observe a linear gradient of base pair probabilities. In sequence 2, the gradient is broken, and instead a central probability value, around 0.5, becomes very frequent. This represents the uniform base pair probability value newly acquired across the whole right arm.


```

AUCAGUUCUAGCAGGAGCUGUACUCAGAGACUCGGGAAUUUUUCCCGGAAUUUUACCCGGGUUUUUACGU
..(((((((.....)))))).....(((((((((((.....(((.....))).....)))))))))).....
AUCGGUUCAGCAGGAACUGUACUCGGGUCUCGGGAAACCCUCCCGGGUUUUACCCGGGUUUUUACGU
..(((((((.....)))))).....(((((((.....(((((((.....)))))))))))))).....
.      *      *      *      *      *      -      00-      *-

```

Figure 6: Seq1 and consensus (seq2) sequence and structure. Mutations are shown in red. The bottom annotation distinguishes mutations in three groups. 1) Those that entail a base pair loss (-). 2) A base pair prevails, even if it's identity is not (*). 3) Mutations in positions with no base pair in either of the structures (0).