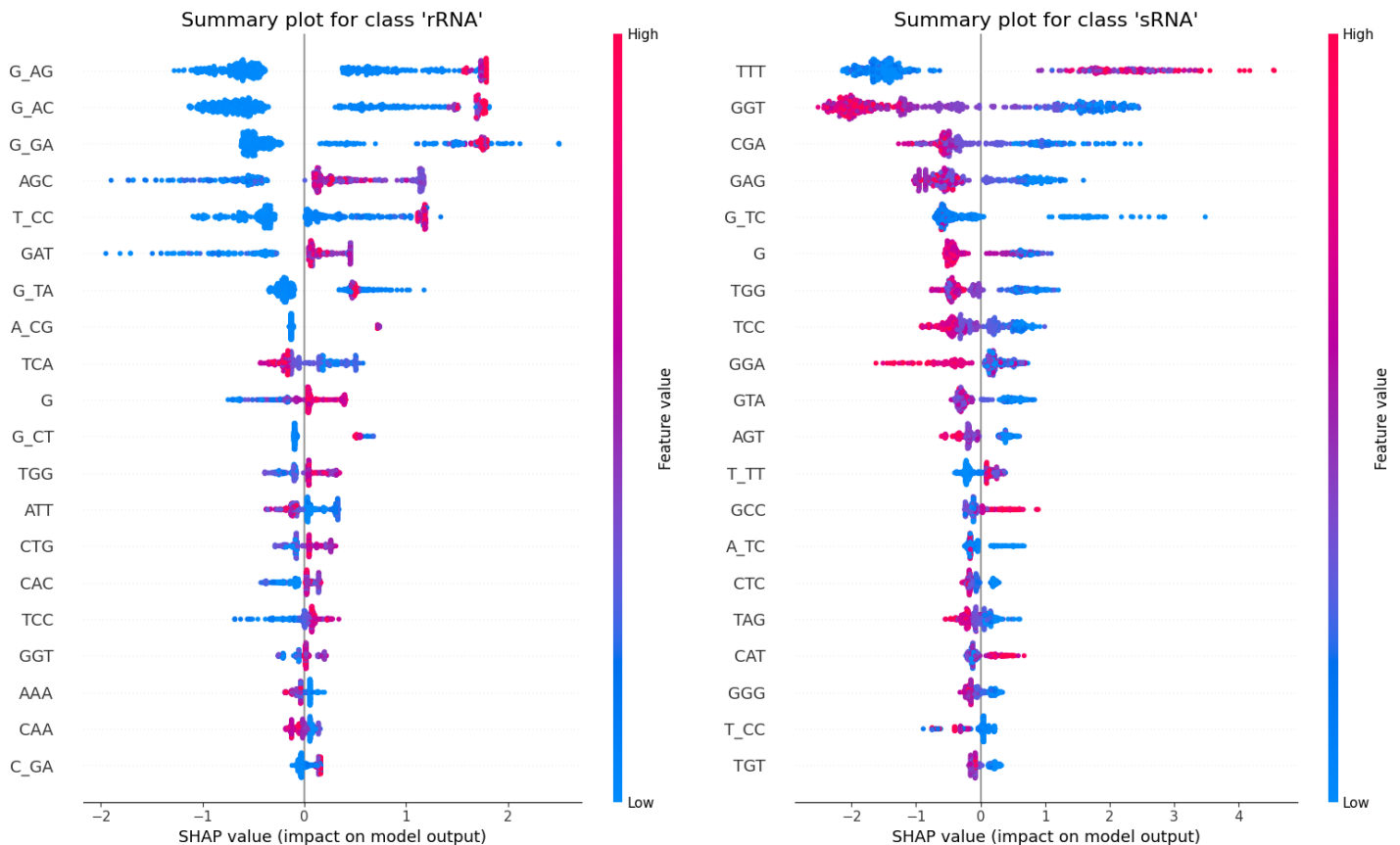
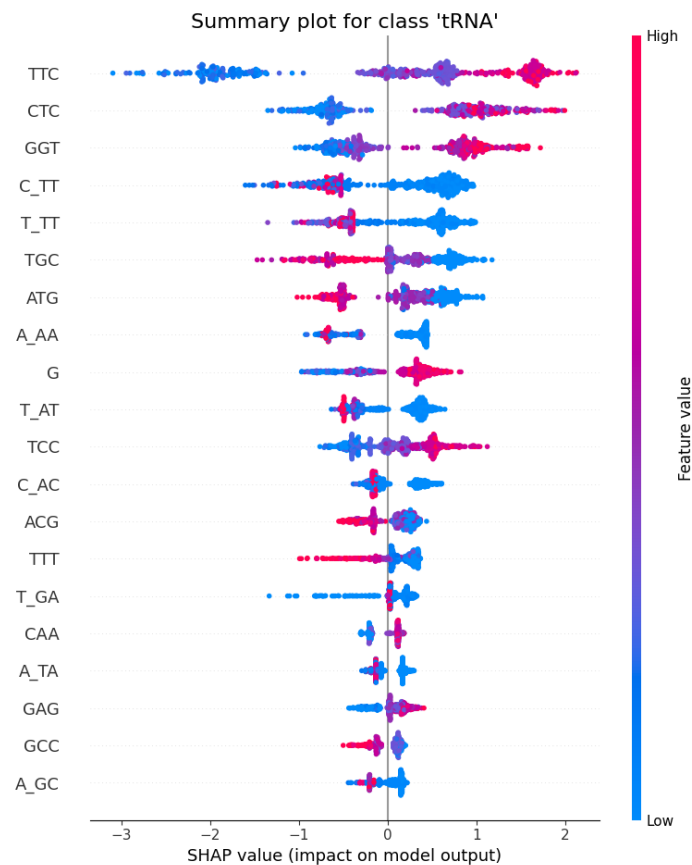


Model Interpretability Report (Multiclass)

This report sustains the idea of being able to interpret and explain how and why the chosen model is classifying each entry as it is. All of Interpretability are based in the SHAP method, in which calculates what's the importance level for each feature in the classification process. It uses Shapley Values, a Game Theory concept, as a descriptive metric to create an hierarquical structure between the features.

Summary Plots





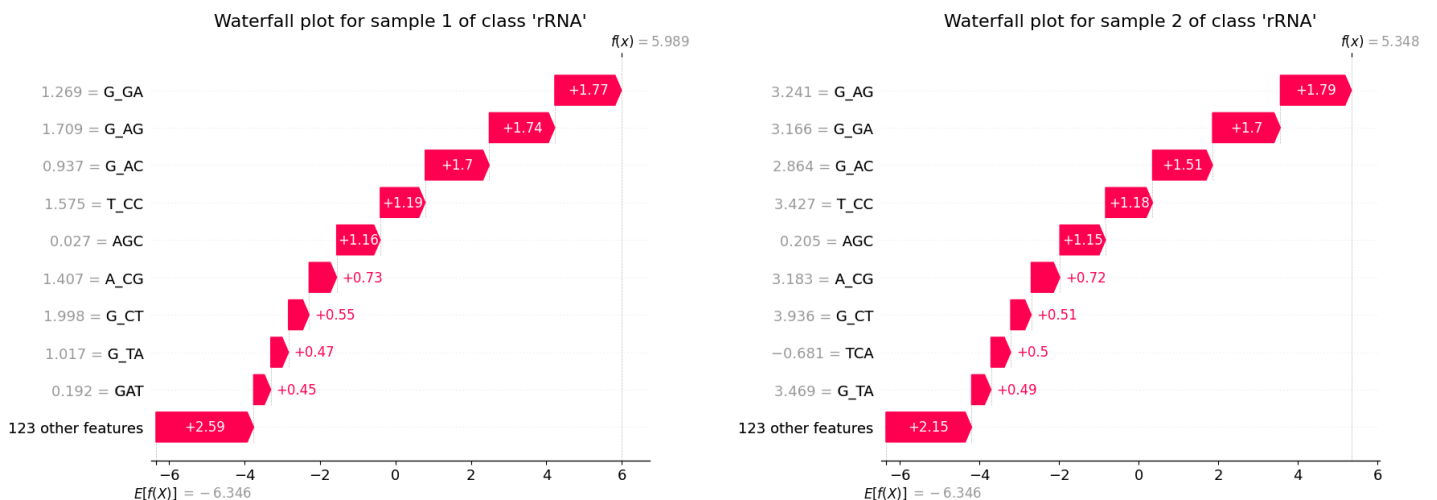
The above plot is called Summary plot and it shows, for each class, how low/high values of each feature contributed for the classification with that class. The features are ranked from most descriptive to least descriptive. This plot is a summarization of all the entries in the test set.

An impact with positive SHAP value means that a high (red dots), medium (purple dots) or low (blue dots) feature value contributes positively of an entry to be classified with that class. The inverse happens with negative SHAP values.

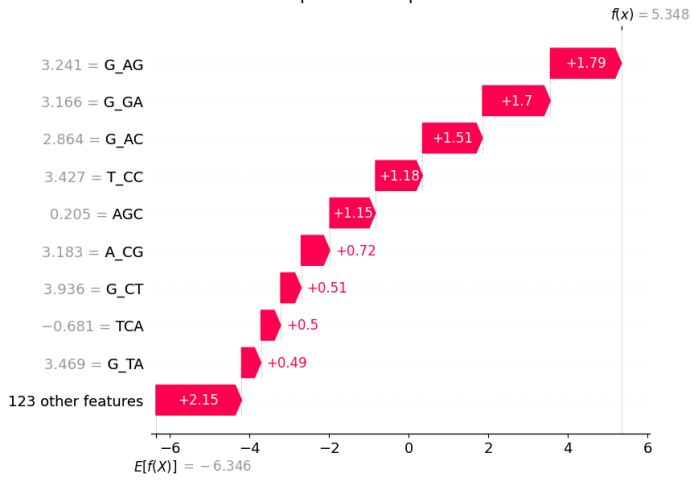
Waterfall Plots

A Waterfall plot shows, for some entries, how each of the features contributed for it to be classified with its classified class. In this case, 3 samples for each class were chosen randomly to be analyzed.

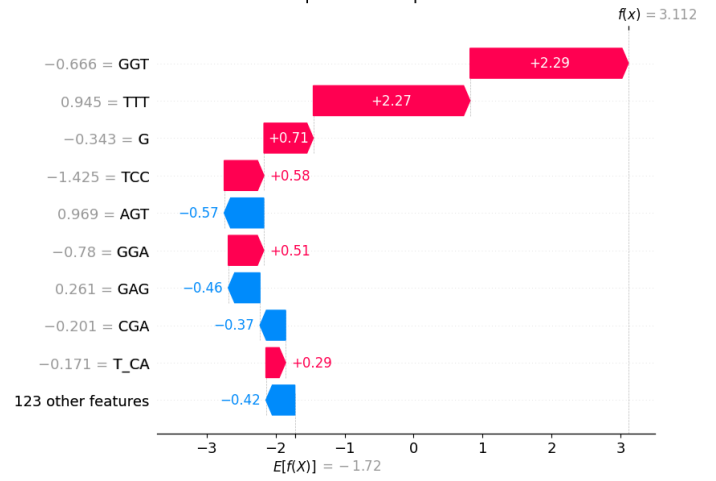
From a base expected value, $E[f(x)]$, each feature contributes positively or negatively towards the entry's given class. At the end, when all the contributions are summed with $E[f(x)]$, we get the final value of $f(x)$ which led to the classification result.



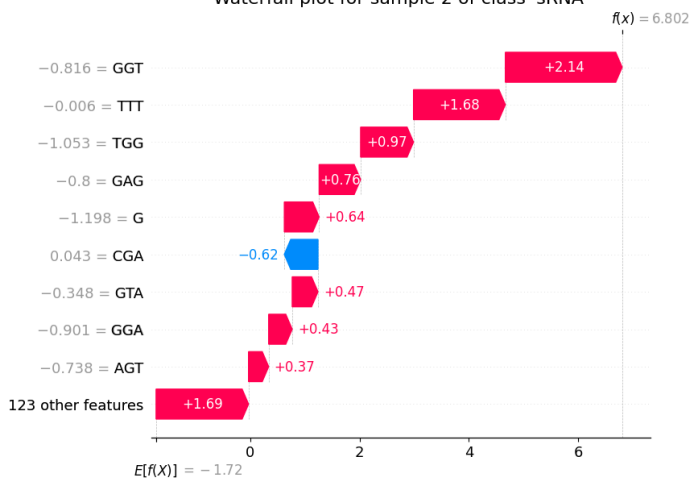
Waterfall plot for sample 3 of class 'rRNA'



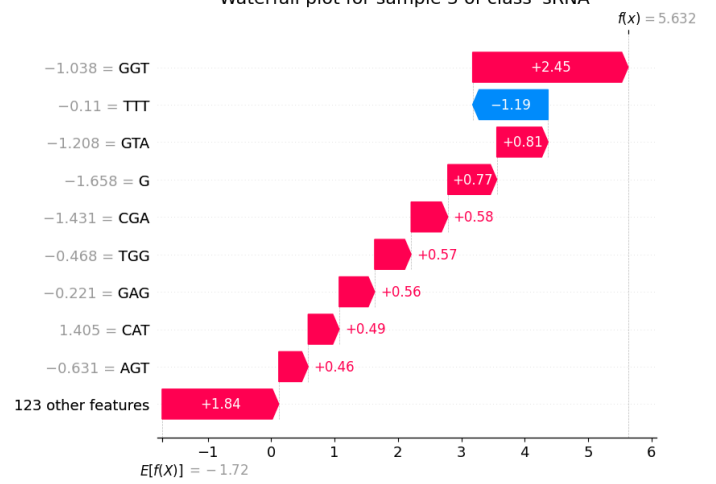
Waterfall plot for sample 1 of class 'sRNA'



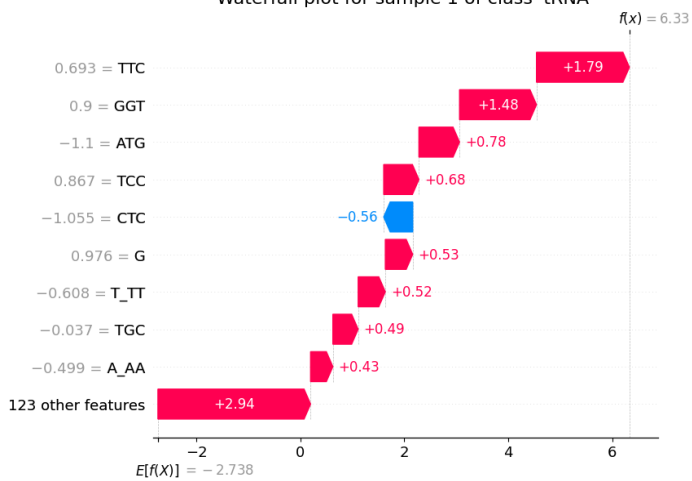
Waterfall plot for sample 2 of class 'sRNA'



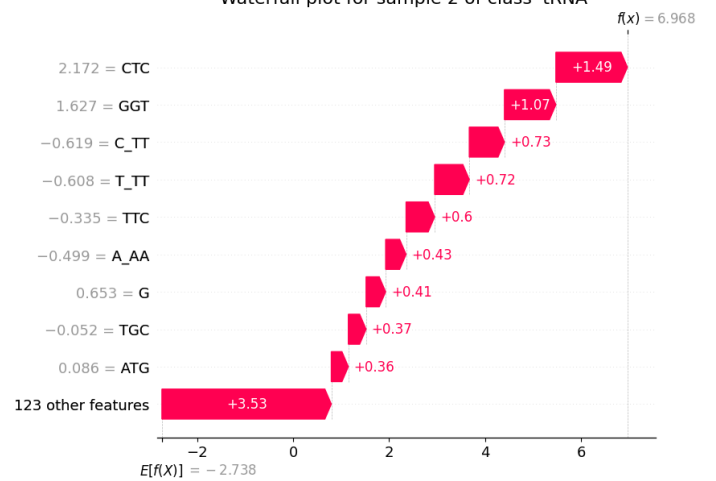
Waterfall plot for sample 3 of class 'sRNA'



Waterfall plot for sample 1 of class 'tRNA'



Waterfall plot for sample 2 of class 'tRNA'



Waterfall plot for sample 3 of class 'tRNA'

