# CLASSIFICATION, PREDICTION AND ANALYSIS OF TYPE VI SECRETED EFFECTOR PROTEINS

{ Rishika Sen, Losiana Nayak & Rajat K. De }
MACHINE INTELLIGENCE UNIT, INDIAN STATISTICAL INSTITUTE, KOLKATA

## 1. WHAT ARE SECRETION SYSTEMS?

1. Pathogens have numerous ways of transporting protein between locations via secretion apparatus [1].
2. Bacterial secretion apparatuses can be divided into different types, based on their structures, functions, and specificity.
3. Secretion systems are used by bacterial pathogens to manipulate the host and establish a replicative niche [2], often ending in infection.

## 2. WHAT IS THE NEED OF CLASSIFICATION?

1. Effector proteins of bacteria infect their hosts by secretion systems (SS).
2. Effector proteins of the T6SS of many species are yet to be discovered.
3. No signature pattern have yet been discovered to differentiate T6 effector proteins from the non-effectors ones.
4. Such signature patterns will speed the discovery of T6 effector proteins in many gram-negative bacteria.

## 3. METHODOLOGY

Features of nucleotide sequences are 4 mono-nucleotide, 16 dinucleotide and 64 trinucleotide frequencies. The protein feature matrix has been made up of 20 single amino acid frequencies, 400 dipeptide frequencies and 11 physico-chemical property based features. We have applied Sequential Forward Selection (SFS) strategy on the feature set of the experimentally verified T6 effector proteins to rank our features according to their significance. According to these ranks, we have created sets of 8, 16, 24, 32 and 40 most significant features. We have compared the performance of the different classifiers with each of these sets and the set with maximum accuracy and stability has been used for prediction of effector proteins in *Vibrio cholerae* and *Yersinia pestis*.
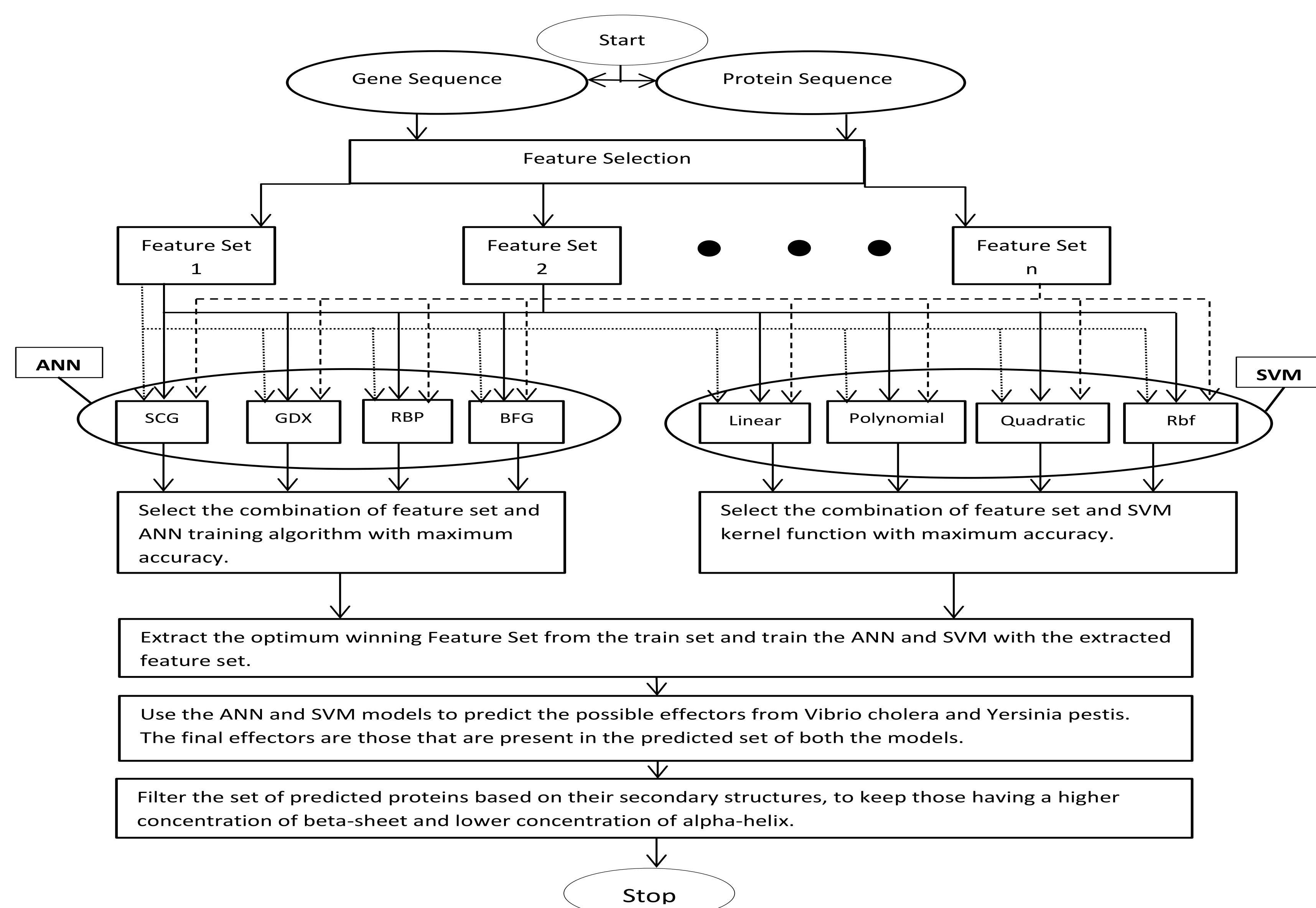


**Figure 1:** Methodology pipeline depicting the flow of the procedure

## 4. CLASSIFIER RESULT

The performance of the training algorithms is as follows:

| Class | SCG | RP | BFG | GDX |
|---|---|---|---|---|
| ANN | 90.02 | 90.74 | 88.53 | 80.21 |

| Class | Linear | Polynomial | Quadratic | RBF |
|---|---|---|---|---|
| SVM | 93.3 | 95.40 | 93.5 | 91.6 |

**Table 1:** Summary of classifier performance

The ANN trained with SCG and SVM with linear kernel have predicted T6 effector proteins with maximum stability. Out of **2267** proteins of *Vibrio cholerae*, **611** proteins have been predicted by both ANN and SVM to be putative T6 effector proteins. For *Yersinia pestis*, out of **2736** proteins, **1058** proteins have been predicted to be possible T6 effector proteins.

## 5. SECONDARY STRUCTURE BASED FILTERING

Table 2 highlights secondary structure based findings of the predicted effector proteins. The ML based predicted effector proteins have been filtered based on secondary structure of proteins of *Vibrio cholerae*, **115** proteins have been found with higher concentration of sheets than helices. Likewise, **193** proteins of *Yersinia pestis* have been found with higher concentration of sheets than helices, as was observed in the experimentally verified T6 effector proteins. These **115** and **193** proteins have a higher chance of being T6 effector protein, since they have a high structural similarity with the experimentally verified effector protein list.

| Class | Coil (%) | Helix (%) | Sheet (%) |
|---|---|---|---|
| Effector | 47.59 | 19.22 | 33.17 |
| Non-Effector | 60.16 | 29.55 | 10.28 |

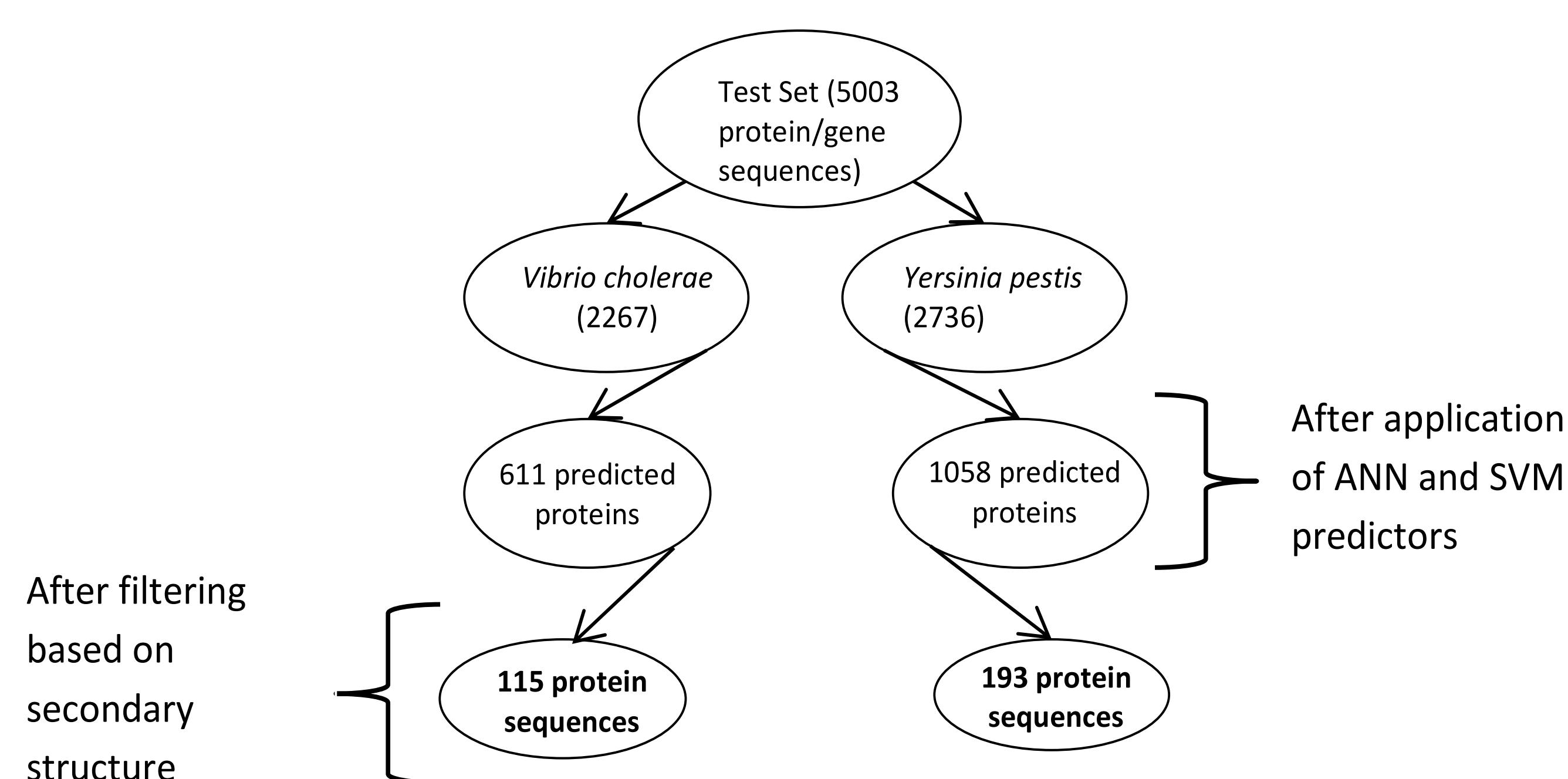**Table 2:** Composition of secondary structure components

## 6. SUMMARY



**Figure 2:** Detailed numerical and species-wise information of the filtering pipeline leading to a more potent set of effector proteins.

We have predicted a set of 611 T6 effector proteins out of 2736 proteins from *Yersinia pestis*. We have also predicted a set of **1058** T6 effector proteins out of **2267** proteins of *Vibrio cholerae*. Among them, **193** proteins of *Yersinia pestis* and **115** proteins of *Vibrio cholerae* have been found with higher percentage of sheets and lower percentage of helices than the non-effector proteins.

These findings have strengthened our prediction, as higher percentage of sheets and lower percentage of helices is a feature often found in other types of effector proteins. The classifier designed here can be further extended to predict effector proteins in other gram-negative (pathogenic) bacteria.

Availability of more experimentally validated T6 effector proteins will generate a better understanding of their structure and functionality, which in turn will strengthen our prediction pipeline.

## 7. REFERENCES

[1] Sen et.al. A review on host–pathogen interactions: classification and prediction. *European Journal of Clinical Microbiology & Infectious Diseases*, 35(10):1581–1599, 2016.

[2] Costa et. al. Secretion systems in gram-negative bacteria: structural and mechanistic insights. *Nature Reviews Microbiology*, 13(6):343–359, 2015.

## THANK YOU!