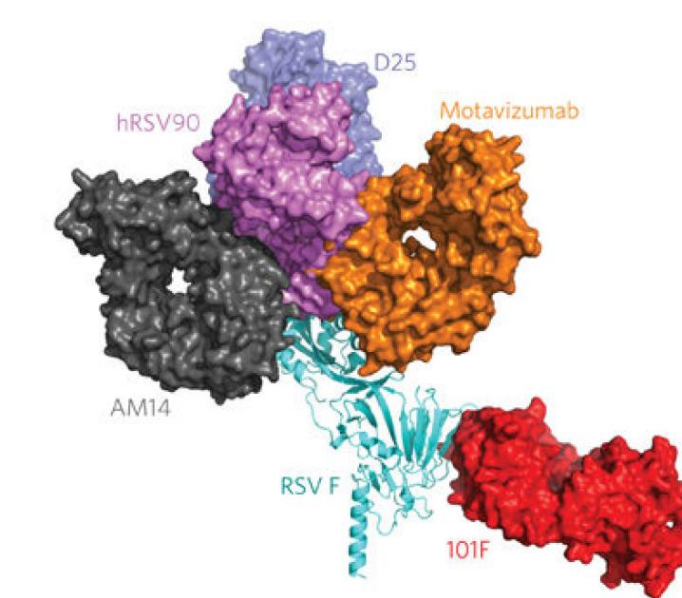


## Introduction &amp; Motivation

Respiratory syncytial virus (RSV) is the primary cause for the lower respiratory infection. To design the vaccine, template-free *de novo* methods is powerful to deal with structurally complex epitope. However, not all designed sequences are necessarily good candidates. Consequently, another preliminary screening process must be performed. Normally, the screening process requires intensive experiments (e.g., Next-Generation Sequencing experiment) [1-2]. This motivates one to design a computational screening process. Thus, this project aims to search for a set of feature and model to approximate a screening function that can determine the quality of a *de novo* designed sequence.

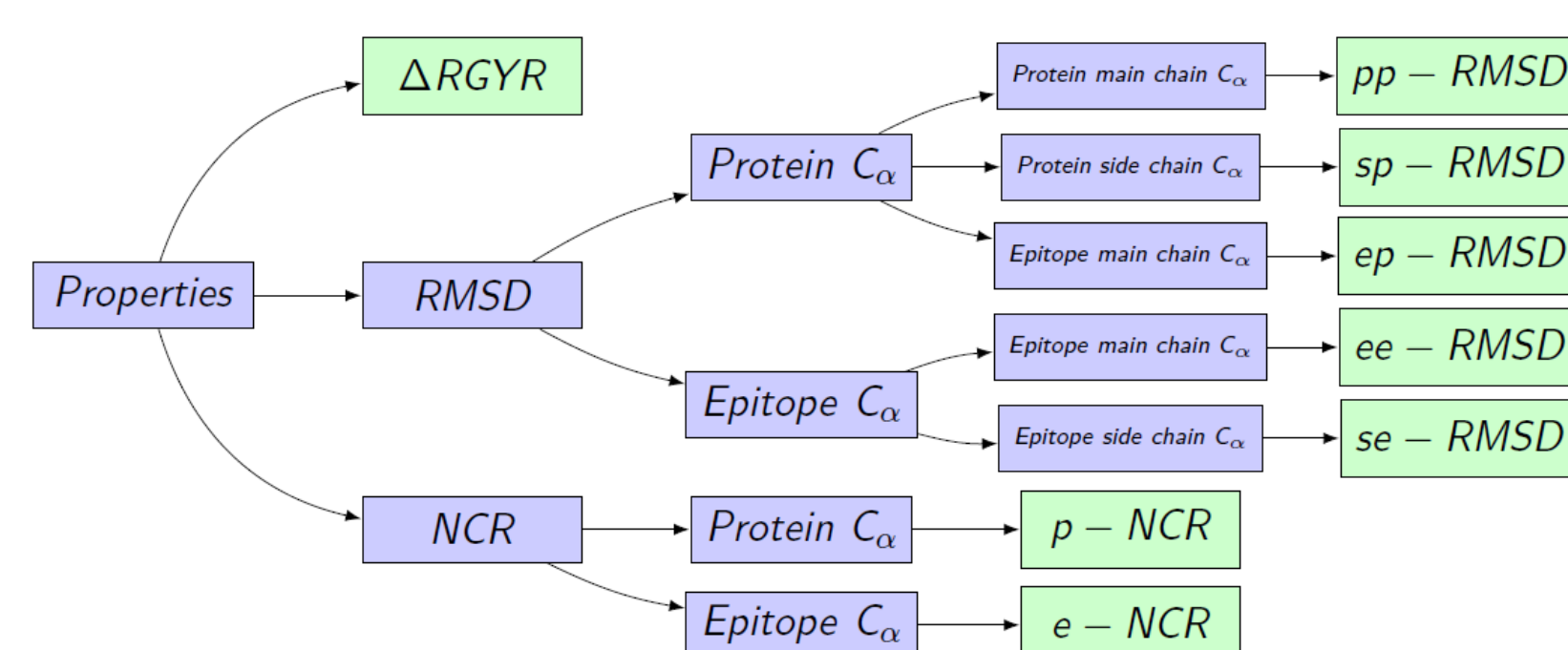


## Trajectory Method

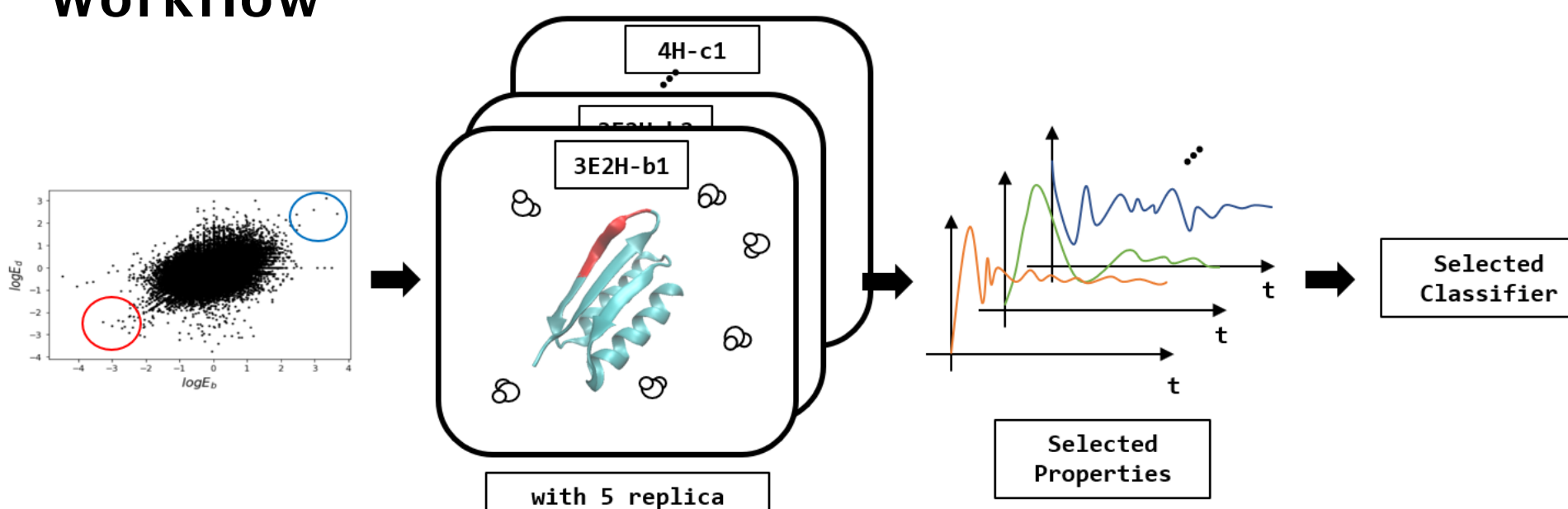
## Methodology

The goal of the trajectory method is to fit the classifier using shape-related ensemble averaged properties obtained from simulation trajectories. Specifically, molecular dynamics (MD) is applied to simulate a *de novo* designed sequence under physiological conditions. A small number of *de novo* designed sequences were first selected, then corresponding structures were generated and simulated by MD. Afterward, shape-related ensemble averaged properties were calculated based on MD trajectory. Finally, the classifiers (Random forest, support vector machine (SVM) and logistic regression classifiers) were trained by these properties.

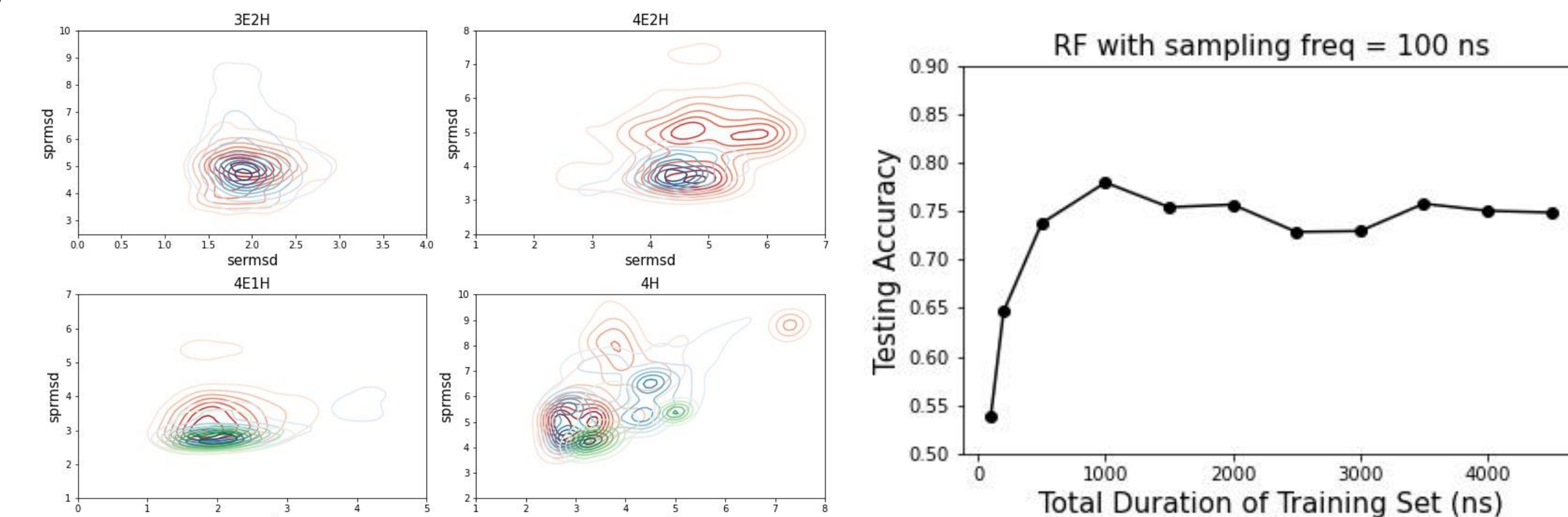
Target Antibody	101F			D25
Design	3E2H	4E1H	4E2H	4H
# Best (g)	4	3	4	3
# Worst (b)	4	3	2	3
# Experiment (c)	0	1	0	1



## Workflow



## Results

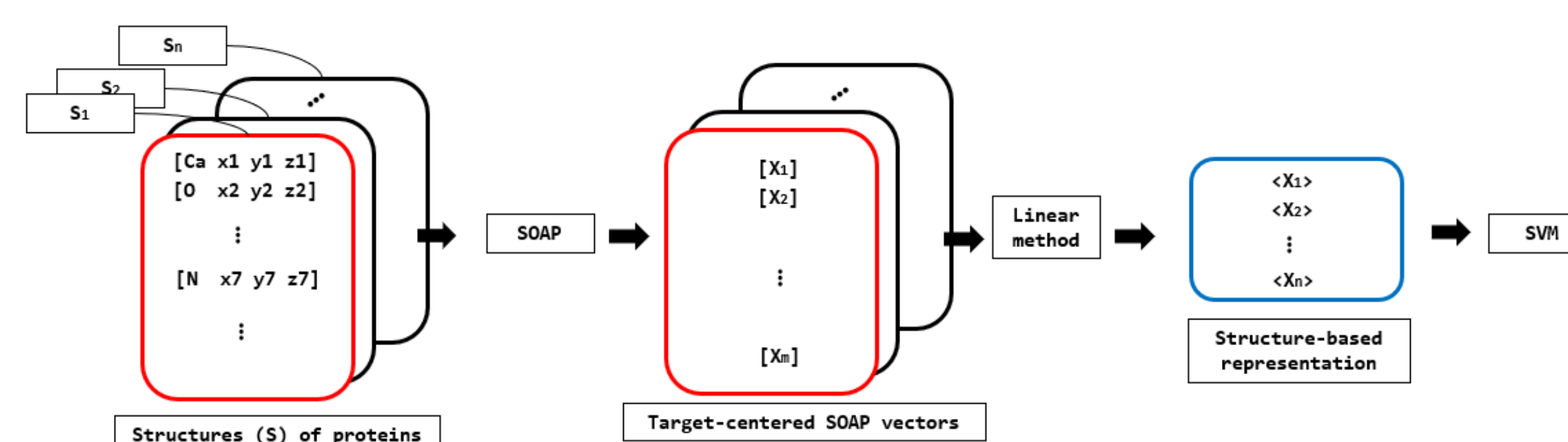


## Structure Method

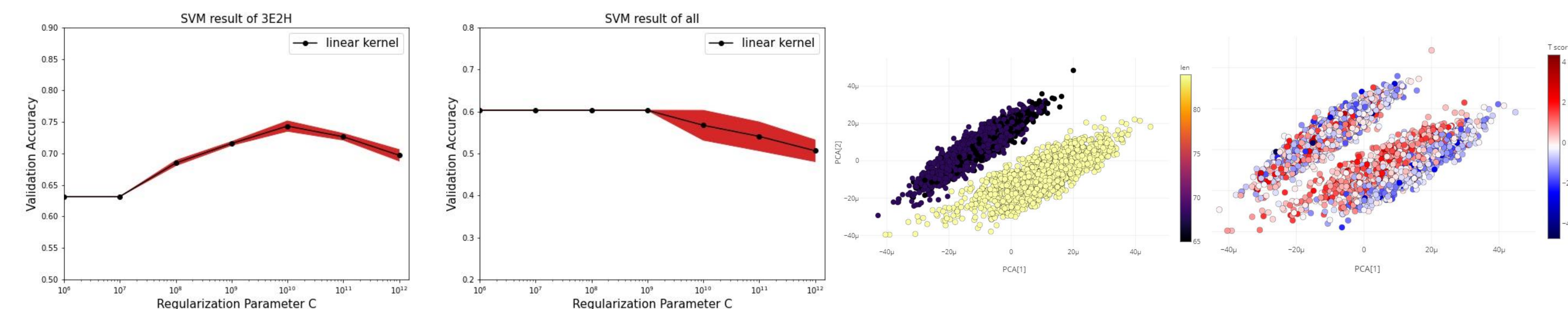
## Methodology

Since shape-related ensemble averaged properties are insufficient to make a satisfying prediction, additional factors should be considered. Therefore, the structure method used SOAP vectors as features. First, all structures (in pdb file format) were originally in Cartesian coordinates. These structures were transformed into SOAP vectors. Subsequently, a linear method was applied to these vectors to obtain a structure-based representation. Finally, these structure-based representations were used as features to train an SVM.

## Workflow



## Results

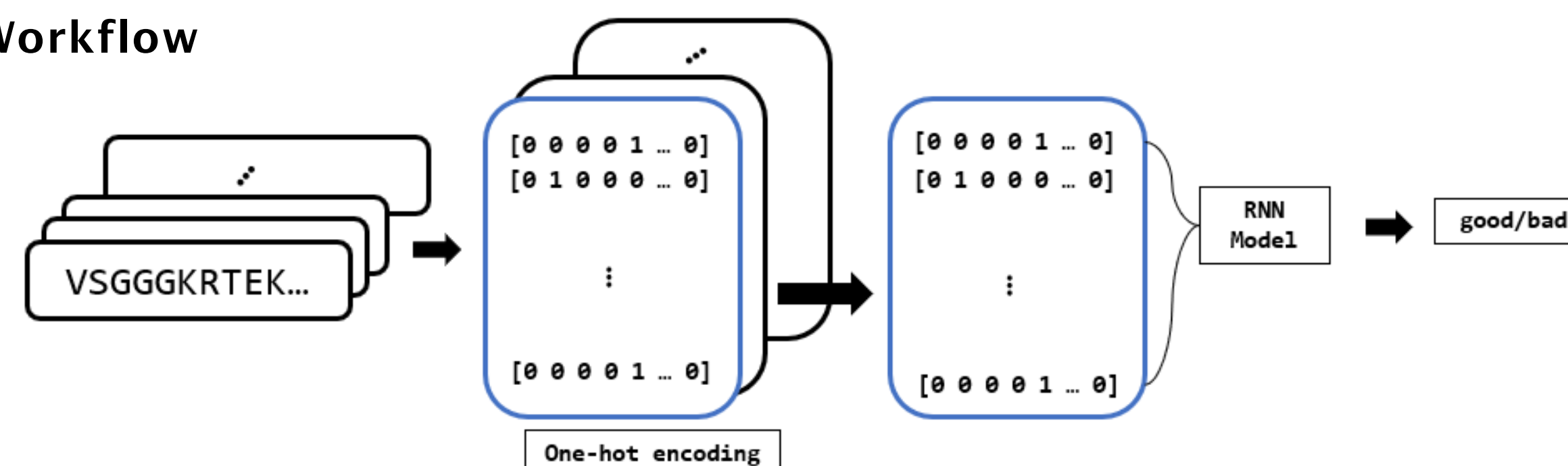


## Sequence Method

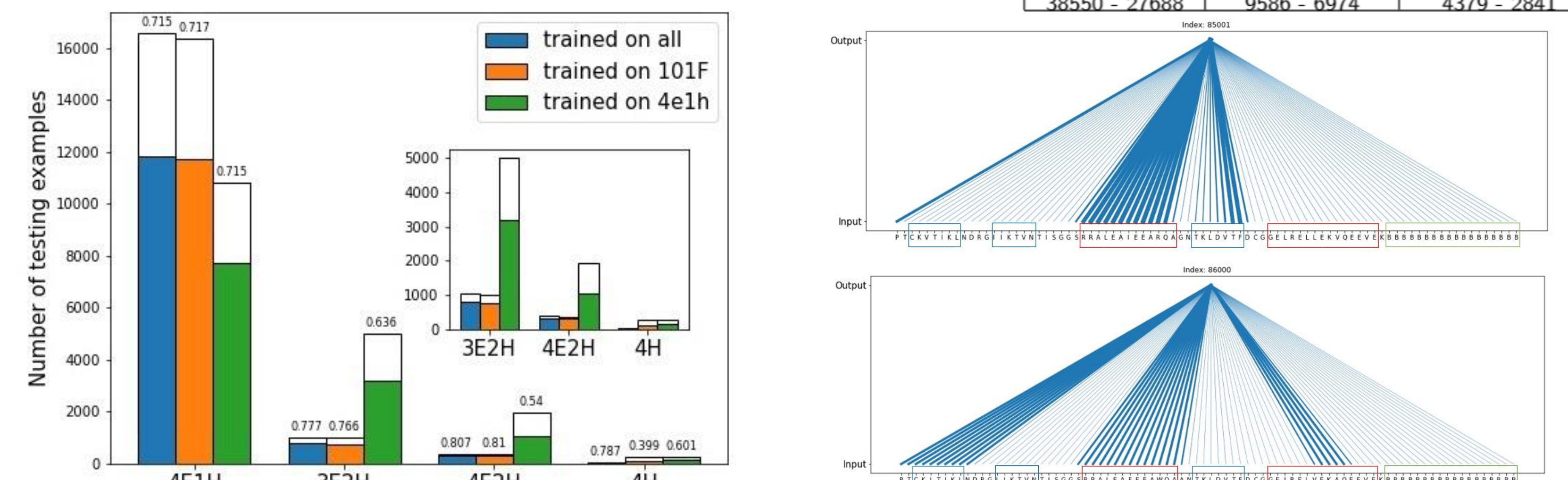
## Methodology

The structure method was only valid when all structures have a single design. In addition, it is itself computationally costly. To find a better model, proteins' sequences were adopted as feature for the sequence method. First, all the sequences were one-hot encoded as input. Next, according to a cross-entropy loss function, the RNN model was trained. Finally, the best model was obtained based on a grid search. The predicted class of given sequences was chosen to be the class that had a larger probability.

## Workflow



## Results



## Conclusion

The trajectory method failed to classify any given sequence. Although the structure method could be useful to some extent, such method was still unrealistic. Finally, the sequence method was able to process 18 times more sequences and twice faster than the structure method. Moreover, the testing accuracy was 0.72 when several designs were included. The classification results showed that the quality of the designed sequence was determined by the sequence's structure regardless of epitope mimicking domain. This disproved previous assumption that sequences with the same epitope-mimicking domain are in the same distribution to the model. Nevertheless, the sequence method could still be helpful because the motivation of this project is to obtain good sequences through computation.

[1] F. Sesterhenn, C. Yang, et al. "De novo protein design enables precise induction of functional antibodies in vivo". In: bioRxiv (2020), p. 685867.

[2] F. Sesterhenn, M. Galloux, et al. "Boosting subdominant neutralizing antibody responses with a computationally designed epitope-focused immunogen". In: PLoS biology 17.2 (2019)