

# Integrating GeoVisualization, Machine Learning, and Deep Learning for Classification and Prediction of Rocky Mountain Spotted Fever in Arizona



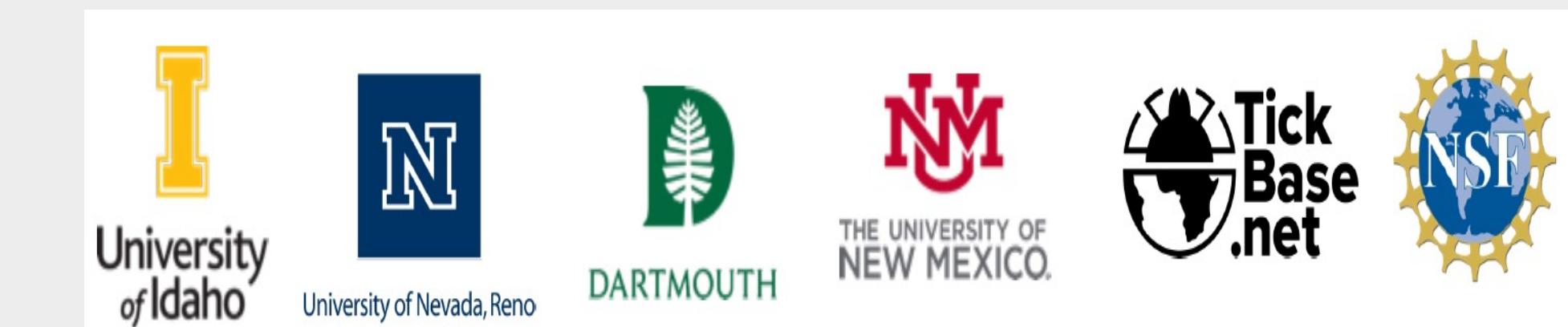
THE UNIVERSITY OF  
NEW MEXICO.

Al Ekram Elahee Hridoy<sup>1,2</sup>, BS; Yan Lin<sup>1,2</sup>, PhD; Zhuoming Liu<sup>2,3</sup>, MS

<sup>1</sup>Department of Geography and Environmental Studies, University of New Mexico

<sup>2</sup>UNM Advanced Spatial Informatics Research and Education(ASPIRE)

<sup>3</sup>Department of Computer Science, University of New Mexico



## PROBLEM STATEMENT

- Rocky Mountain Spotted Fever (RMSF) is a potentially fatal Tick Disease (TD) caused by the bacterium Rickettsia and has disproportionately affected communities in southwest US (e.g., Arizona, and Navajo Nation). To address this problem, we have developed a machine and deep learning classification and prediction system that uses geovisualization techniques to identify high-risk areas for RMSF in Arizona.
- Our system includes a first-of-its-kind climate-based classification model with a Receiver Operating Characteristic Area Under the Curve (ROC AUC) of 0.97, which can accurately predict which counties in Arizona are most likely to have a high incidence of RMSF. Additionally, we have implemented a deep learning method known as Long Short Term Neural Networks(LSTM) to predict the possible occurrence of RMSF in each county.
- Our approach to predicting the spread of RMSF will enable public health officials to take timely preventive measures and allocate resources to the areas at greatest risk. Our machine learning system has the potential to significantly reduce the incidence of this potentially fatal disease in the southwestern United States.

## OBJECTIVES

- To visualize RMSF disease statistics in a web application.
- To implement climate-based classification system to map most likely to have a high incidence of RMSF.
- To implement deep learning based prediction of RMSF incidence system for each county.

## DATA AND DATA SOURCES

We used yearly time-series data from 2006 to 2021.

- RMSF cases. Data Source: ARIZONA DEPARTMENT OF HEALTH SERVICES
- Land Surface Temperature (Degree Celsius). Data Source: MODIS Land Surface Temperature and Emissivity (MOD11)
- Relative Humidity (In Percentage) and Precipitation (mm, daily total). Data Source: GRIDMET: University of Idaho Gridded Surface Meteorological Dataset
- NDVI (Unit Less). Data Source: MODIS Terra Daily NDVI
- CDC Social Vulnerability index (SVC) and Total Population. Data Source: CDC/ ATSDR Social Vulnerability Index

## GEOVISUALIZATION WEB INTERFACE

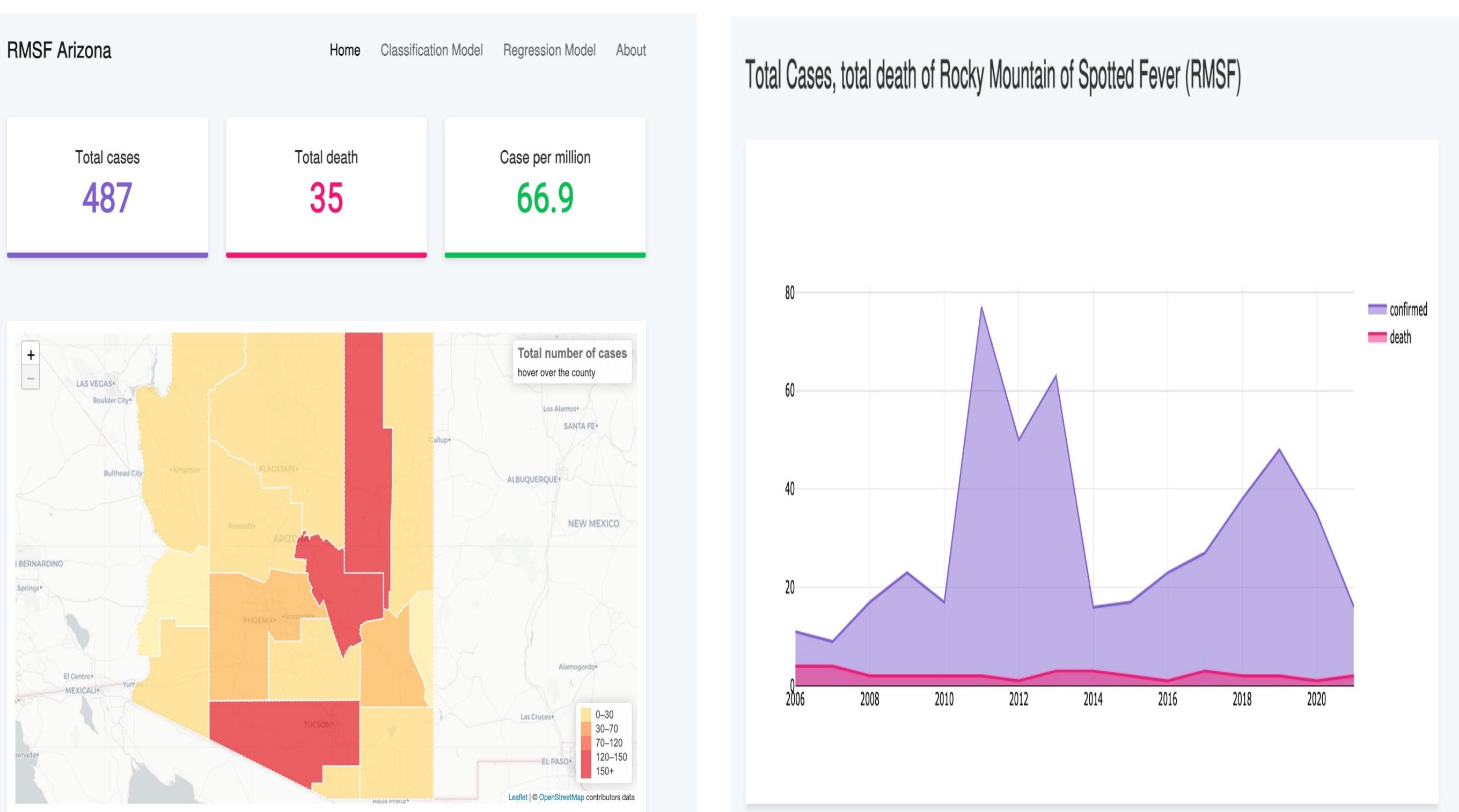


Fig. 1: Geovisualization web interface

## WEBSITE

<https://alhridoy.github.io/RMSF-dashboard>

## METHODS

- We experimented with various machine learning methods to solve our classification problem and ultimately developed a Random Forest algorithm pipeline. Prior to hyperparameter tuning, our pipeline achieved an accuracy of 0.95 with an ROC AUC of 0.82. We were able to improve the performance of our pipeline by tuning the hyperparameters and implementing 5 k-fold cross validation, resulting in an ROC AUC of 0.974. We addressed the issue of class imbalance between high incidence and low incidence using the RandomOverSampler technique.
- For prediction of RMSF, we used a deep learning technique, Long Short Term Neural Networks(LSTMs). The architecture of the model consists of two LSTM layers, each with 64 units and a ReLU activation function, followed by a single Dense layer with one unit, and the model is compiled with the mean absolute error (MAE) loss function and the Adam optimizer .
- For classification we followed the CDC's definition of "high incidence" as a county with more than 10 cases per 100,000 population, the data was binned for classification models by assigning a value of 1 to counties with an incidence rate greater than or equal to 10/100000, and 0 to those with a lower incidence rate.

## RESULTS AND DISCUSSION

- After performing GridSearch, we decided to use the Random Forest algorithm for our classification model, which achieved an ROC AUC of 0.974, a recall of 1.0, and a precision of 0.42. We then tested our model with 2021 data and found that Gila County had a 94.06% chance of having a high incidence of RMSF, while Navajo County had a 36.05% chance. We also identified that Land Surface Temperature (LST) and population were the most important predictor variables.
- LSTM results suggest that test Root Mean Square Error (RMSE) values ranging from 0.006 to 9.055.

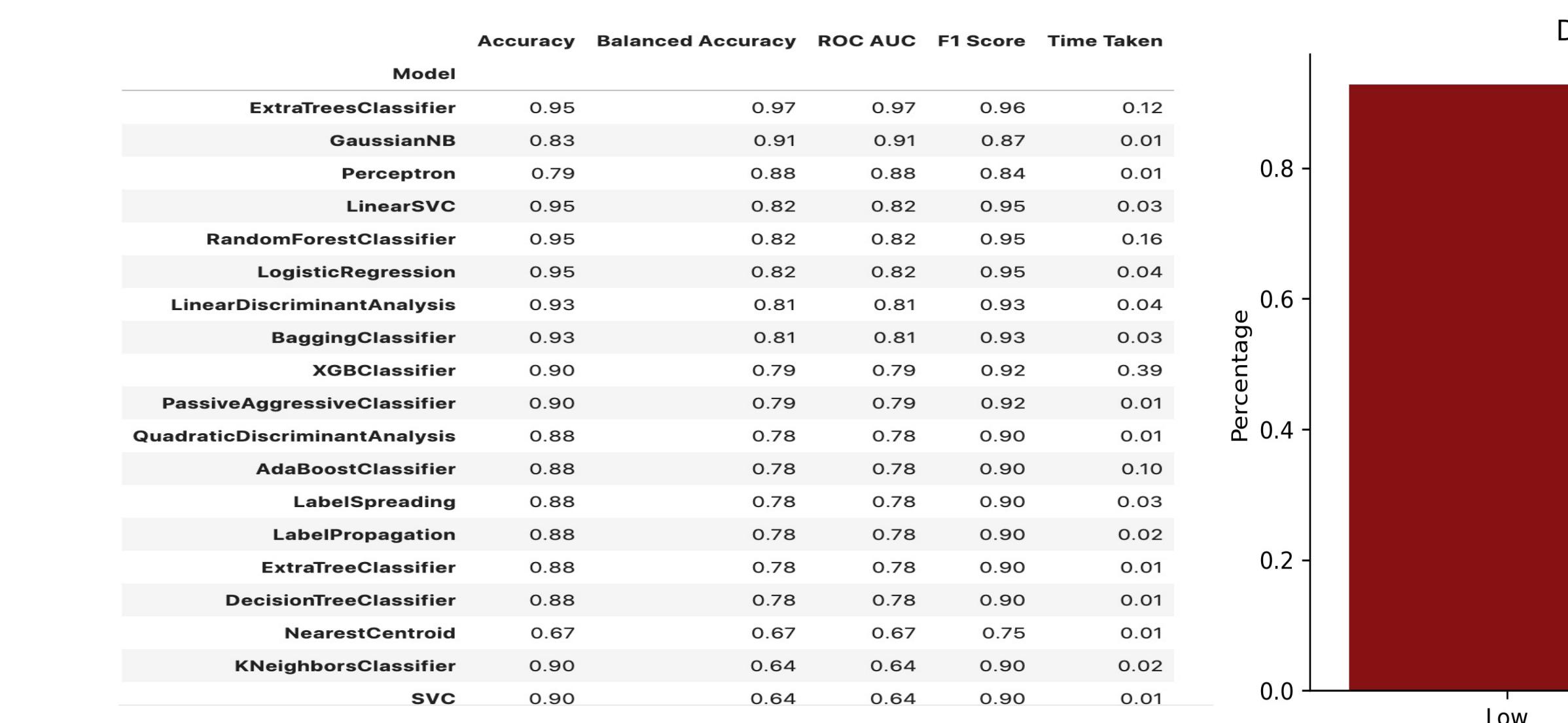


Fig. 2: Machine Learning algorithms for classification of RMSF incidence

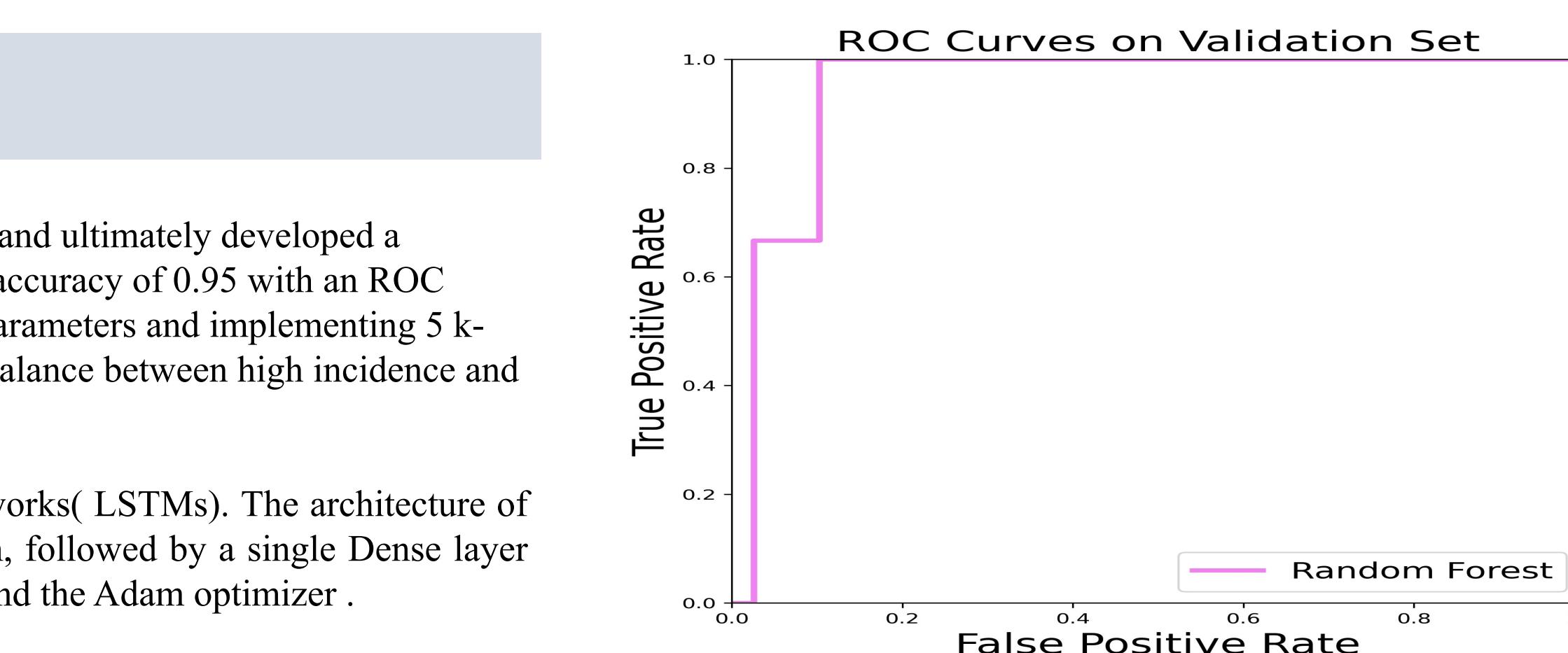


Fig 5: ROC Curve on Validation Set

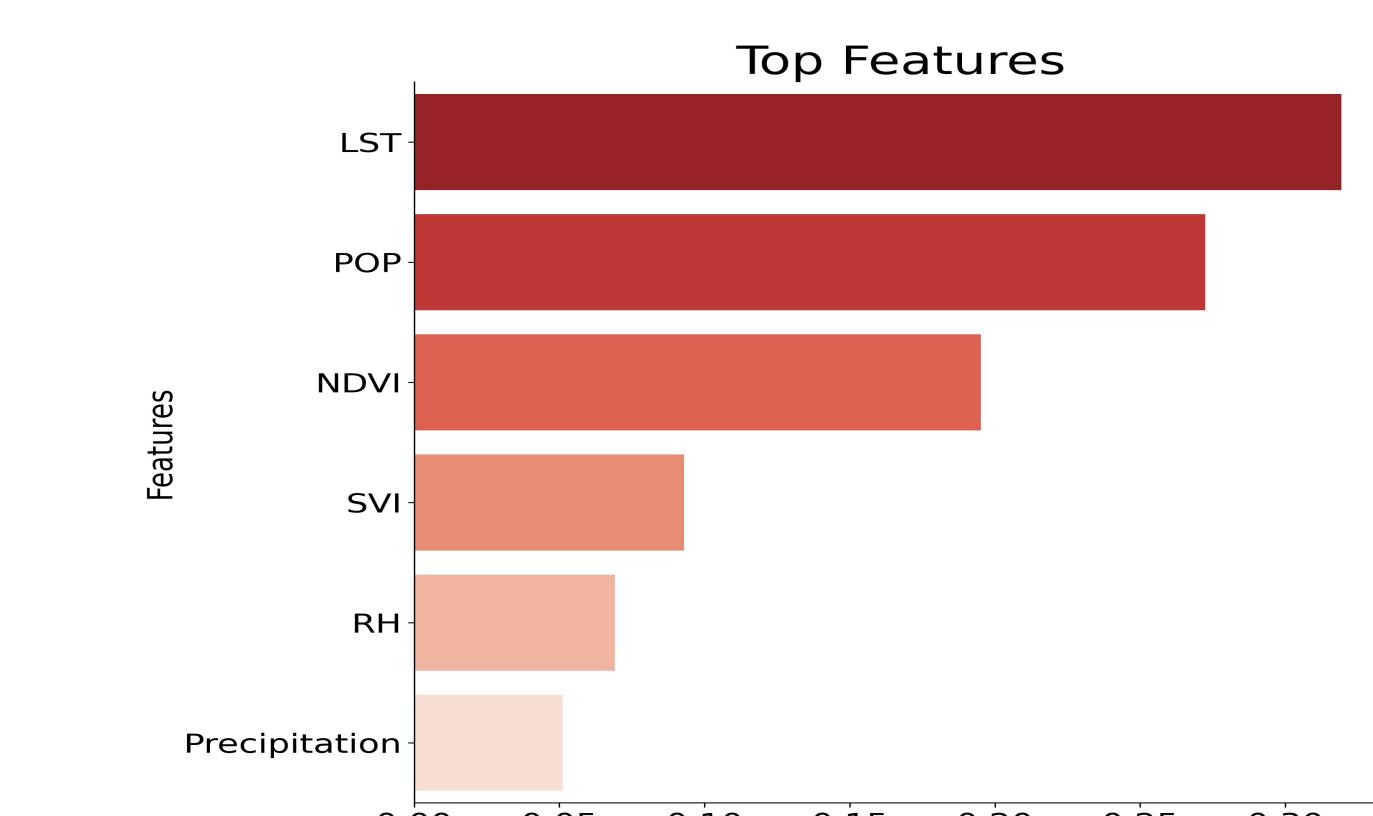


Fig 7: Feature importance of RF

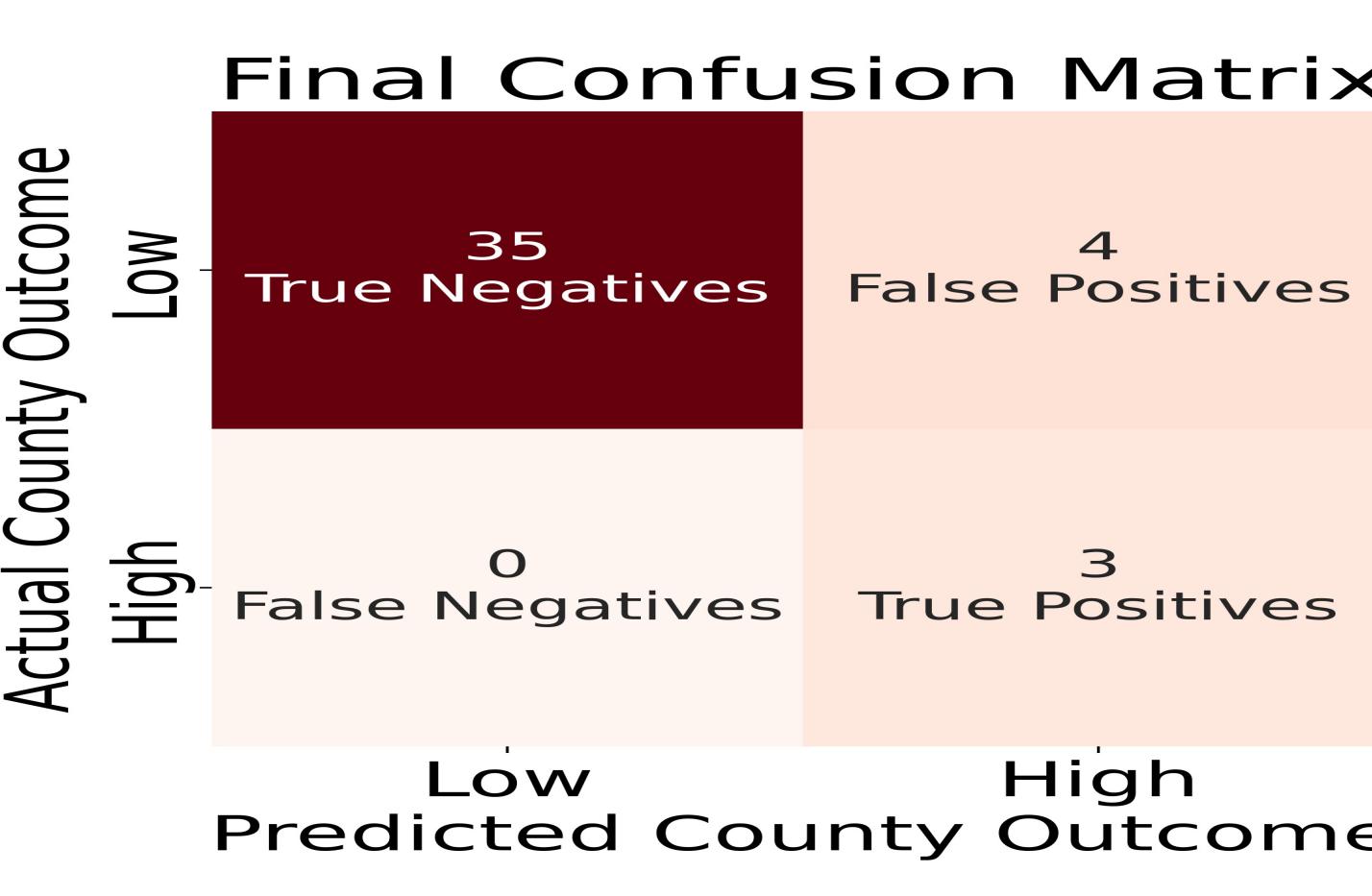


Fig 6: Final Confusion Matrix

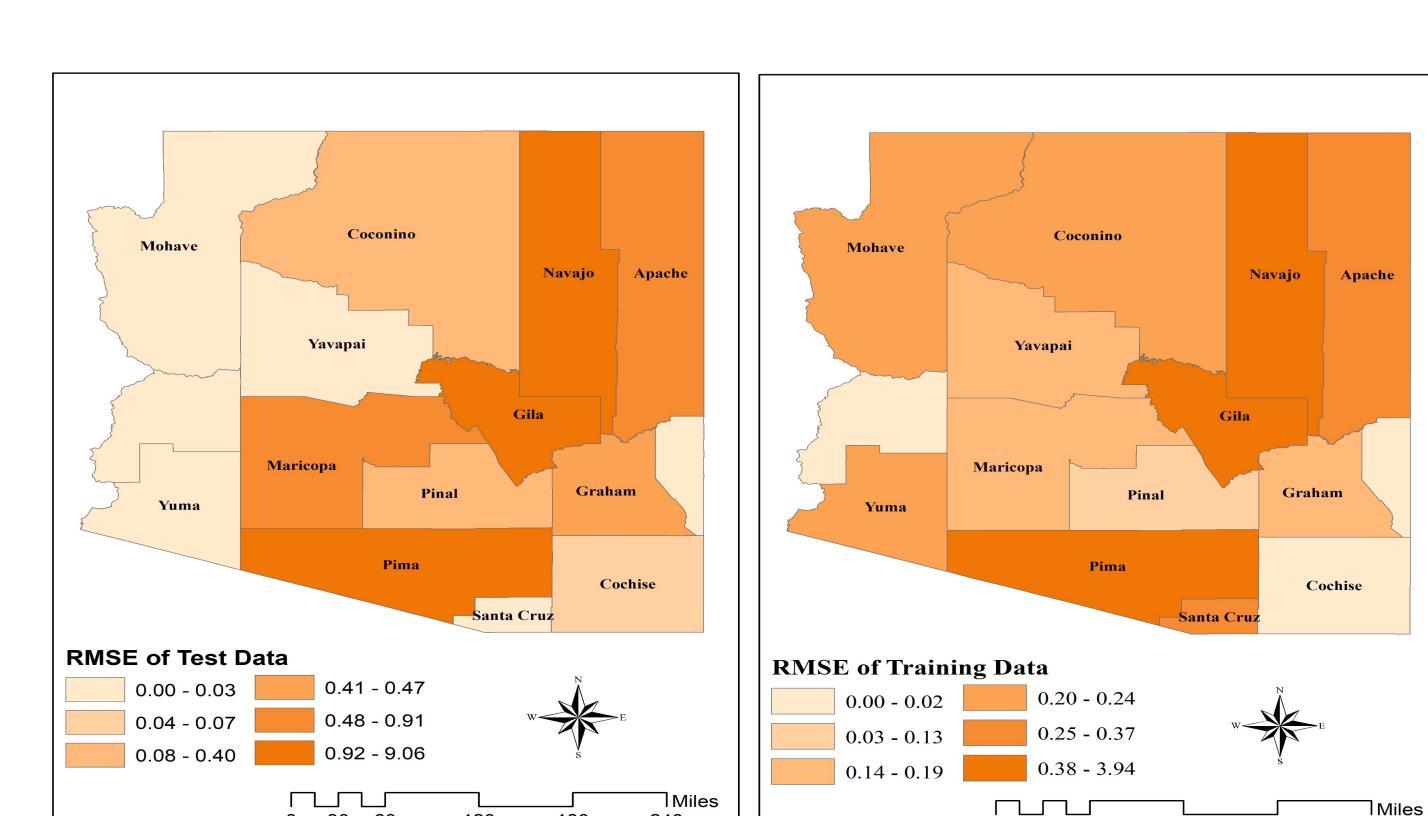


Fig. 8 : RMSE of test and training data

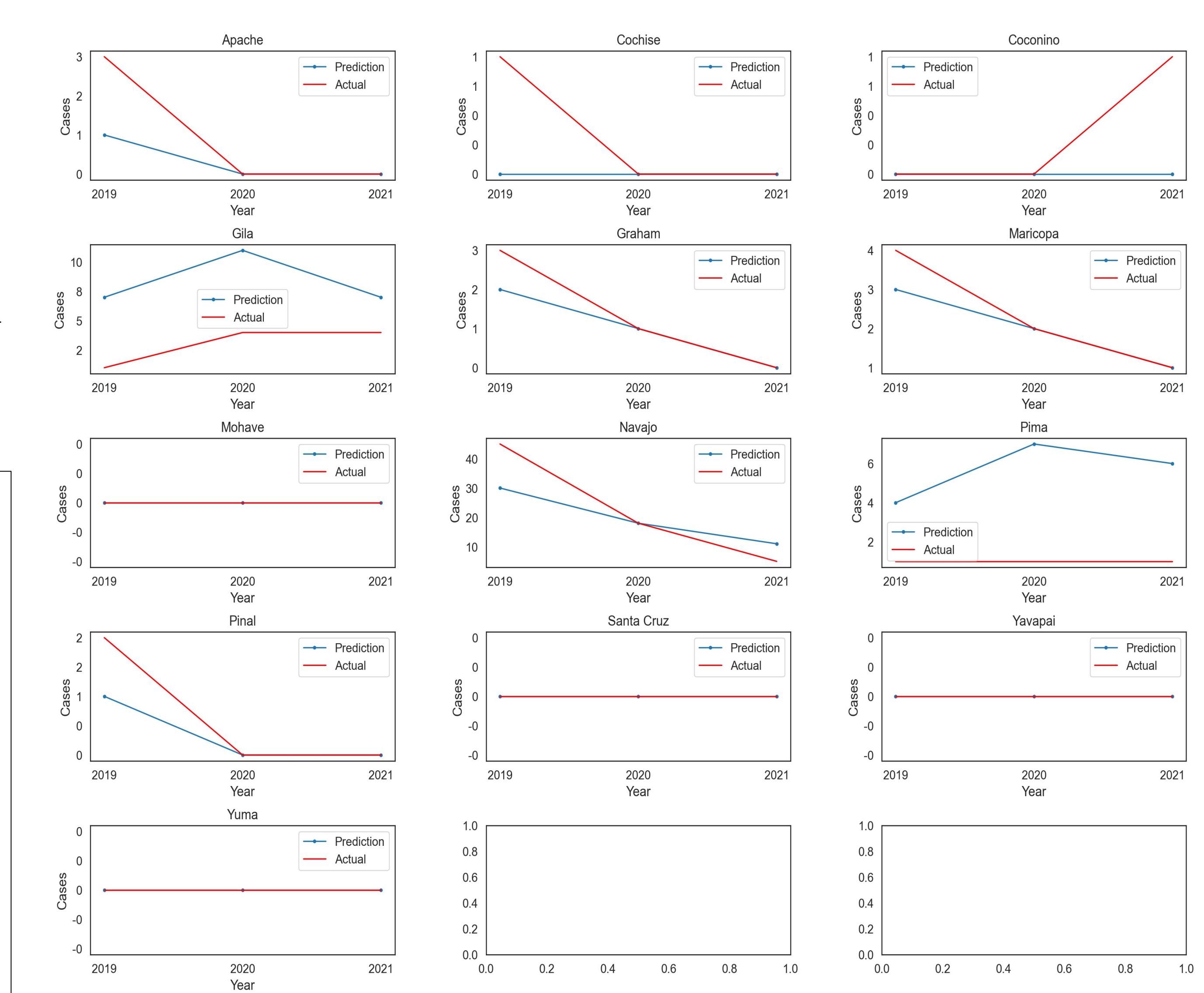


Fig. 9 : LSTM prediction on 3 years of test data 2019-2021

## LIMITATION AND FUTURE WORK

- The data is satellite derived ;thus, temporal and spatial resolution of data is a limitation of this study.
- We had county-level RMSF cases. Moreover, granular-level data will facilitate more accurate modelling of the disease dynamics.
- In future work, authors want to integrate more granular RMSF data at zip-code level to improve the prediction model.

## ACKNOWLEDGMENTS

- This work was facilitated with Google Earth Engine and ADHS, we thank to GEE team and ADHS.
- Funding for this project provided by the National Science Foundation, Award Number 201960.
- This work has not been formally reviewed by NSF. The views expressed are solely those of the authors and do not necessarily reflect those of these agencies.