

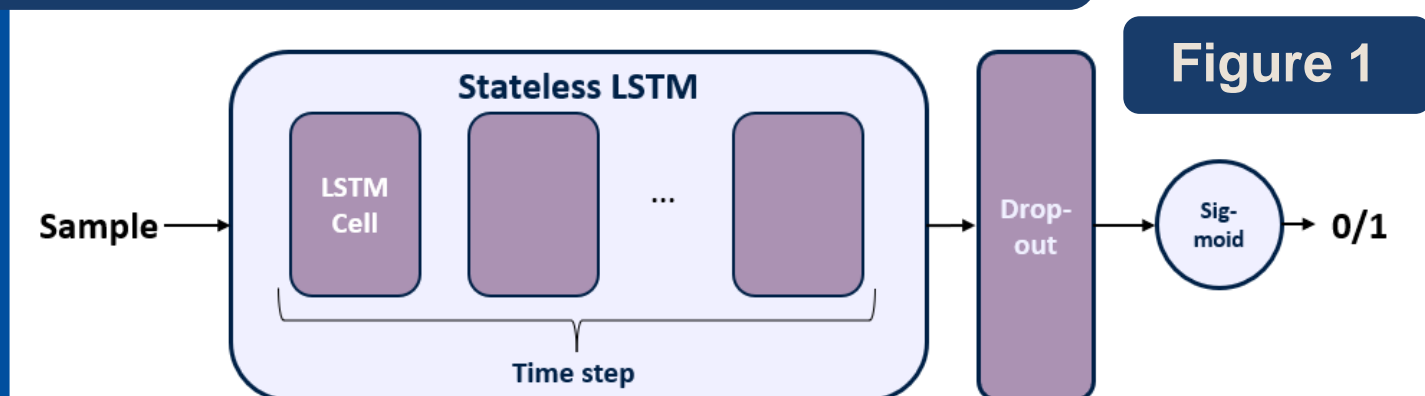
### 1. Introduction

- Dengue is an infectious disease originally endemic in tropical countries such as South East Asia. Recently, Dengue was found in over 188 countries across different continents, putting 3 billion people at risk annually.
- If Dengue progress to Severe Dengue, the mortality rate can be up to 50%, and reduce to 1% if patients get proper treatment.
- Dengue places a heavy burden on the healthcare systems and economy of endemic countries. In Vietnam, every year 2 million people are infected and cost the economy \$95 million annually.

### 2. Objective

Using clinical data collected over 10 years by OUCRU, a machine learning model is hoped to be developed, deployed and aid clinicians in diagnosing Severe Dengue and relieve the strain of overcrowding in hospital

### 3. Machine Learning Model



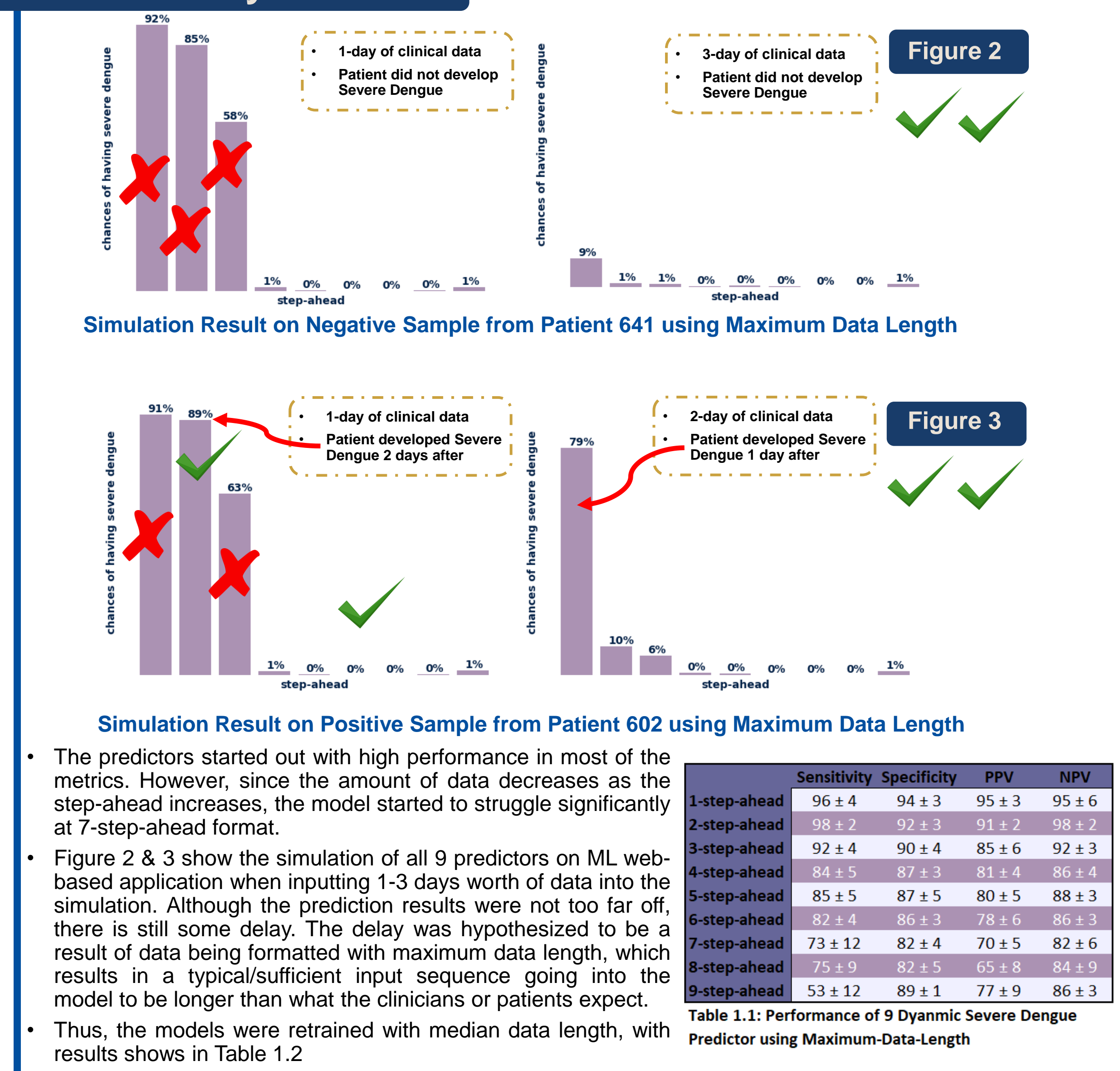
- The first Dengue ML predictive model to use RNN based neural network architecture.
- LSTM is proven to work well with clinical data (\*) through multiple research. It is capable of processing time-series data. And unlike traditional RNN, it is able to decide which information to withhold for long-term and which for short-term.
- Model is optimised with RMSProp algorithm which is good in preventing vanishing and exploding gradient in complex network

### 4. Data Processing

study ID	day from onset	vomiting	body temperature	respiratory rate	haemoglobin	haematocrit	platelet count	bleeding vaginal	bleeding mucosal	abdominal pain	Label
2547	0	0	0	0	0	0	0	0	0	0	1
2547	4	0	0	0.234	41.4	128	0	0	0	0	0
2547	5	0	0	0.234	47.2	132	0	0	0	0	0
2547	6	0	0	0.234	47.5	106	0	0	0	0	0

- Data Imputation:** to simulate the reality that not all medical tests are done on a daily basis, missing values from data is not interpolated but only filled with 0 values. Since 0 does not associated to any meaningful value in any numerical variable, it should not hinder the learning ability of LSTM model (\*\*)
- Data formatting:** Data is formatted to different type known as step-ahead format. 1-step-ahead format means the data is to be used to predict if Severe Dengue will happen in roughly 1 day. Each sequence is then padded to the decided common data-length so the training process would be easier.
- Data Standardisation**

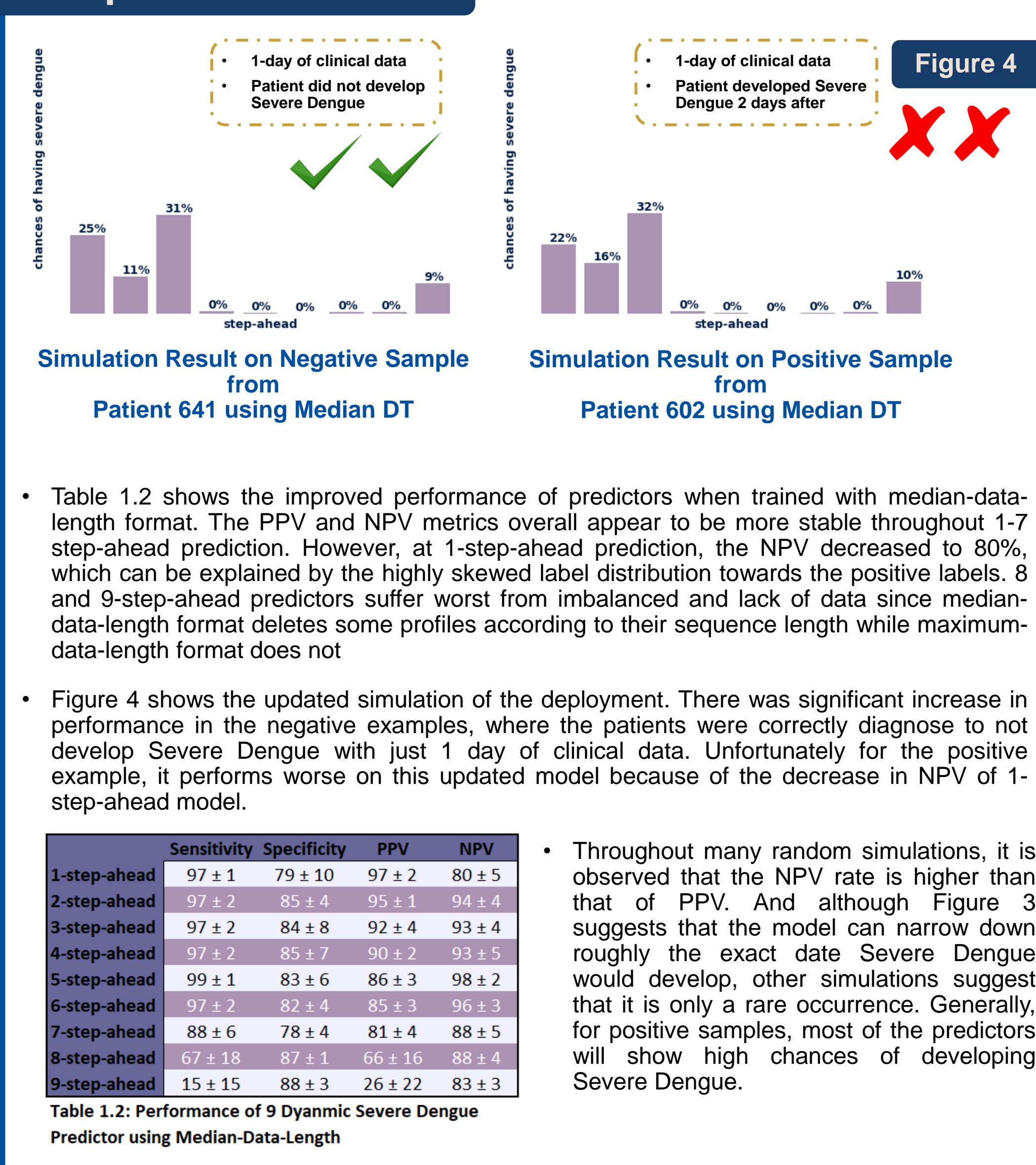
### 5. Preliminary Results



### 7. Discussion

- From the performance results shown in Table 1.1 and 1.2, the predictors shows promising results in classifying Severe Dengue cases. With adequately high PPV and NPV, the deployed model is hoped to help clinicians in clinical management, relieve burden on hospital systems during an outbreak, and giving severely at-risk patients with proper treatment. This result demonstrates the ability of LSTM in learning the dependencies between variables, through times, even when the training data is missing up to 50% data points.
- Figure 2 shows a more desirable performance of the model where for a positive sample, the model was able to narrow down the date that Severe Dengue would develop. Since the current sequences of data are often filled with 0, it is suggested that the lack of data points (missing values) hinders the LSTM ability to expand on its hypotheses. However, during simulation, the models still exhibit performance consistent with test result shown in Table 1.1 and 1.2
- For future work, a model trained on imputed data set would be a good starting point. Other data formats which promote the reality of short input data being provided into the models should also be explored. A smaller set of features will also increase efficiency of the since they would require less amount of clinical test which may be costly and time consuming. Instead of binary classification, a regression model that predicts the day Severe Dengue is worth looking into. In this project, the regression model was attempted. However, date features were either missing a lot of data points or incomprehensibly recorded, which hinder the feasibility of developing and training the regression model.

### 6. Updated Results



(\*) J. Xia, S. Pan, M. Zhu, G. Cai, M. Yan, Q. Su, J. Yan, and G. Ning, "A long shortterm memory ensemble approach for improving the outcome prediction in intensive care unit," Computational and Mathematical Methods in Medicine, vol. 2019, pp. 1-10, 2019.  
 (\*\*) T. Pham, T. Tran, D. Phung, and S. Venkatesh, "Predicting healthcare trajectories from medical records: A deep learning approach," Journal of Biomedical Informatics, vol. 69, pp. 218-229, 2017.

(\*\*) F. Chollet, Deep learning with Python. Manning, 2018  
 NPV: Negative Predictive Value  
 PPV: Positive Predictive Value

