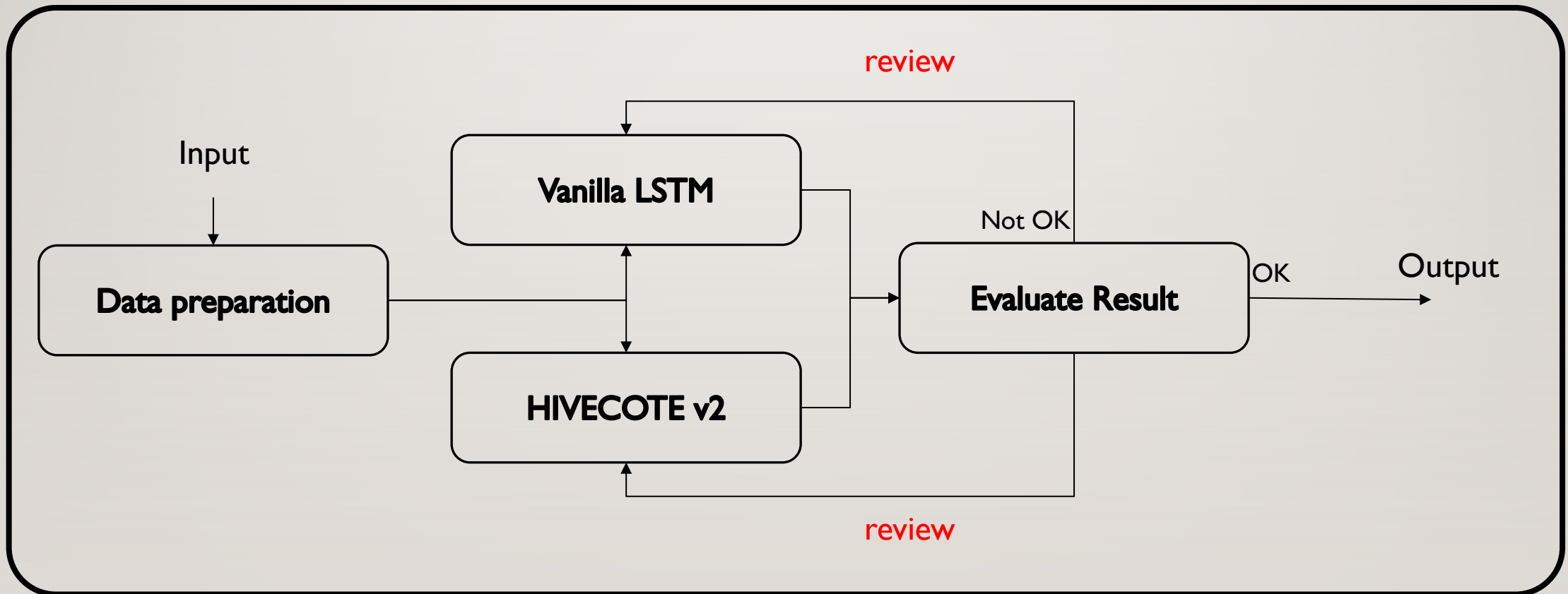# DATA SCIENCE PROJECT

- Zhang
- Zihao Zhang

# CONTENT

- Define the Problem

- Workflow

- Data Preparation

- Modelling

- Further Improvement

# DEFINE THE PROBLEM

- What is the problem?
  - 370 samples belonging to 9 persons, classify 270 unknown samples
  - Each sample is a 7~29 * 12 dimensional time series (12 LPC)
  - Time series multi-class classification

- How to solve the problem?
  - Vanilla LSTM
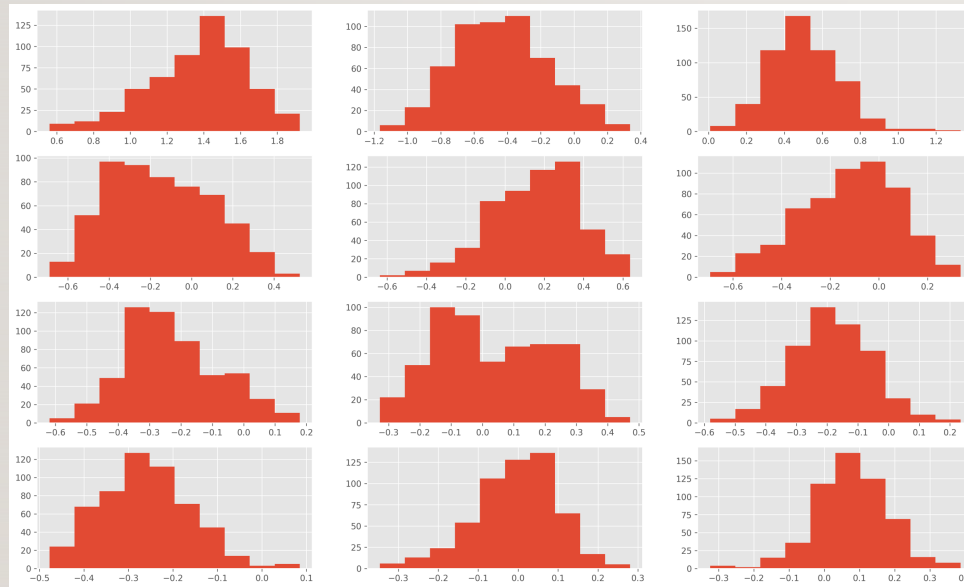  - SOTA model – HIVECOTE v2.0

# WORKFLOW

# METRIC

- **Accuracy** for all classes

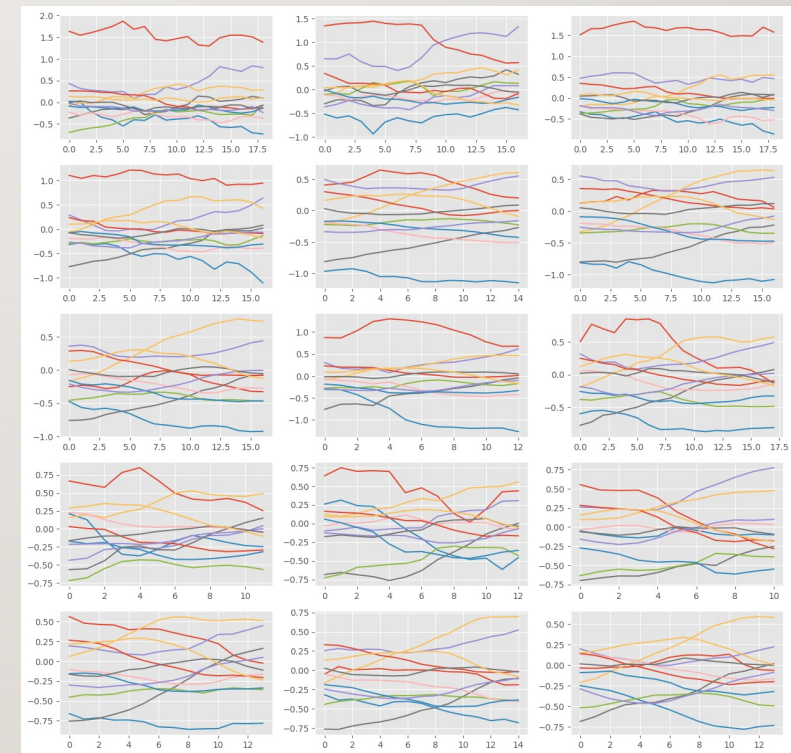- Worst performance user's **accuracy**

# DATA PREPARATION

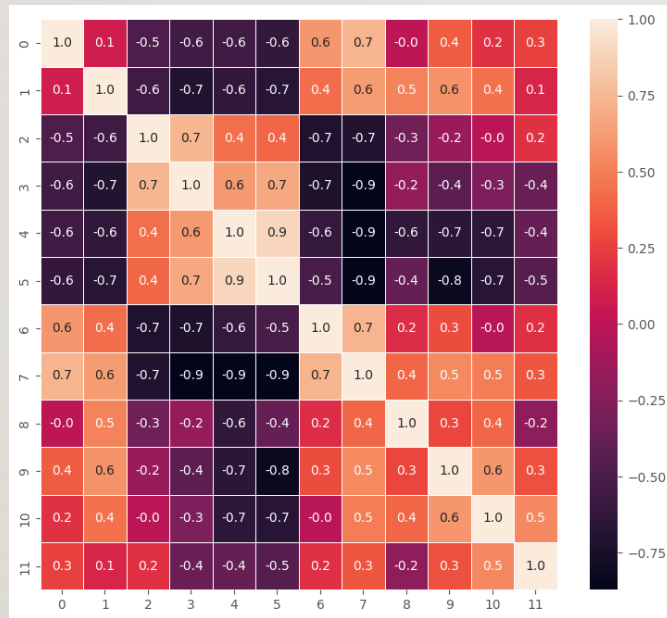- Exploratory Data Analysis and Visualization

12 LPC histogram for person 1

12 LPC line-plot for head 15 samples
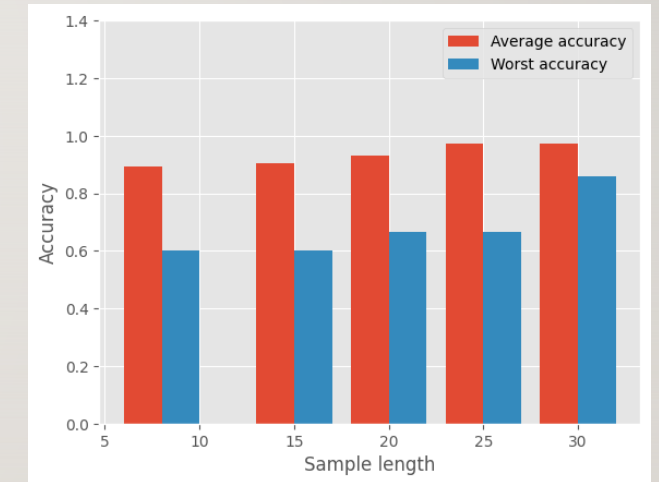
# DATA PREPARATION

- Exploratory Data Analysis and Visualization
  - High autocorrelation

# DATA PREPARATION

- Padding: to make all data the same length
  - Maximum length, padding with last value – keep more information
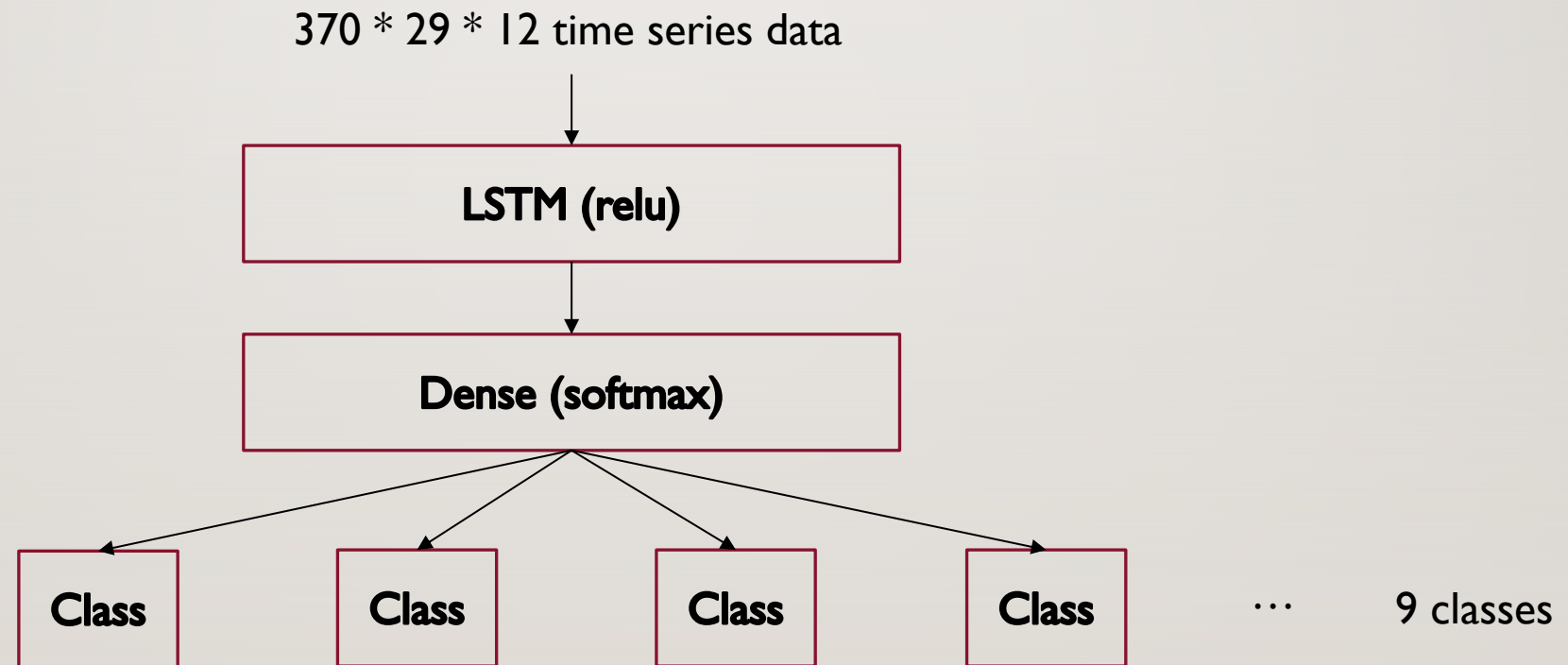  - Average length ×

$$X_i \qquad \begin{matrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \\ x_{i5} \\ x_{i5} \\ x_{i5} \\ x_{i5} \end{matrix} \qquad X_j \qquad \begin{matrix} x_{j1} \\ x_{j2} \\ x_{j3} \\ x_{j4} \\ x_{j5} \\ x_{j6} \\ x_{j7} \\ x_{j7} \end{matrix} \qquad X_k \qquad \begin{matrix} x_{k1} \\ x_{k2} \\ x_{k3} \\ x_{k4} \\ x_{k5} \\ x_{k6} \\ x_{k7} \\ x_{k8} \end{matrix}$$

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - Maximum length

# MODELLING

- Vanilla LSTM

370 * 29 * 12 time series data

↓

| LSTM (relu) |
|---|

↓

| Dense (softmax) |
|---|

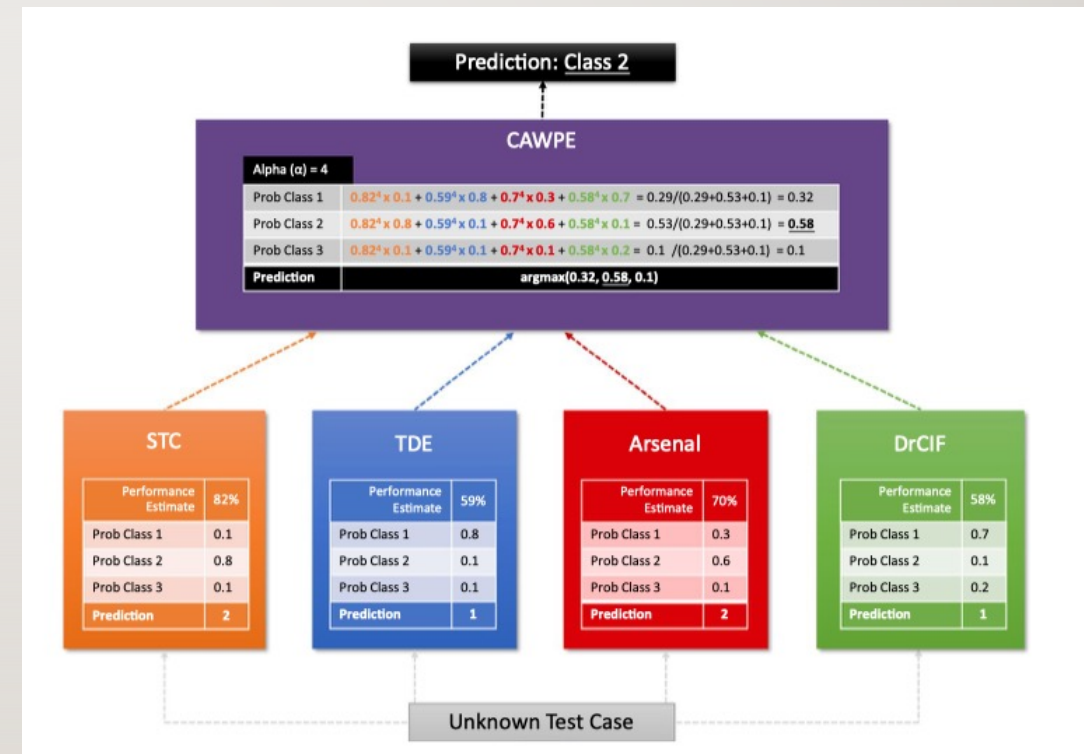| Class | | Class | | Class | | Class | | ... | 9 classes |
|---|---|---|---|---|---|---|---|---|---|

# HIVECOTE V2

- State of the art

- Ensemble model
  - Shapelet Transform Classifier
  - Arsenal
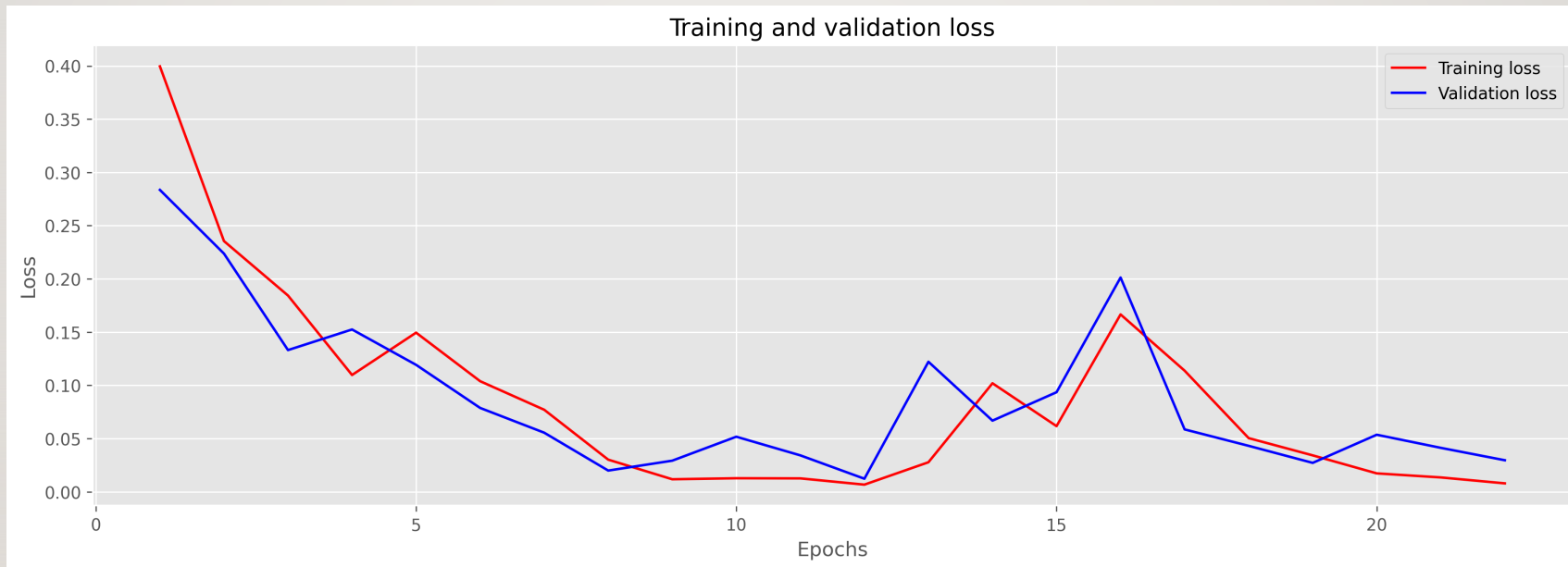  - Dictionary based representation TDE
  - the interval based DrCIF



Picture source: https://arxiv.org/pdf/2104.07551.pdf

# IMPROVE THE MODEL

- Augmentation:
  - Weighted resample augmentation
  - All dataset augmentation
  - One class resample augmentation
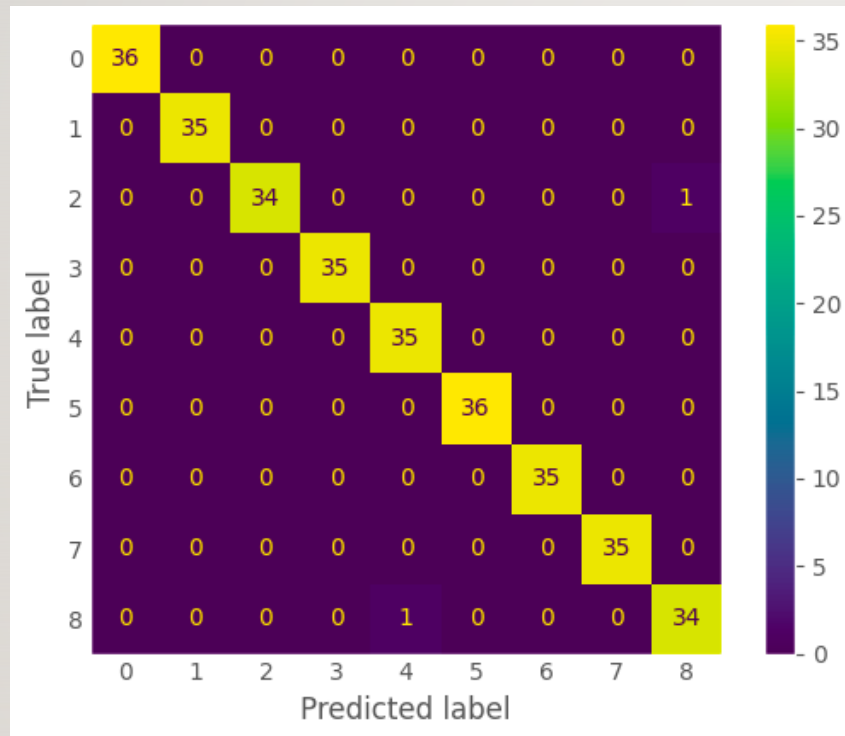
Yield best results

# VANILLA LSTM



BEST validation set performance:
- The classification accuracy is 0.9874
- The worst performing user accuracy is 0.9429

# HIVECOTE V2



Confusion matrix

BEST validation set performance:

- The classification accuracy is 0.9937
- The worst performing user accuracy is 0.9714

# RESULT IMPROVEMENT

- Final result is not robust and worst user prediction accuracy vibrates (best save to csv) add regularization

- More feature engineering

- Gather more data from the last person

- Better data augmentation methods (slicing, warping, jittering, rotation, and their combination)

- Hyperparameters tuning

- HIVECOTE v2 is slow, especially for large dataset (replaced by RocketClassifier)

- Data Leakage problem

- Focal Loss (add BinaryFocalCrossentropy as the evaluation criteria)

- DO NOT DO
    - Include test data with labels (easy to guess) for training ☺

# THANK YOU!

- Zihao Zhang