**Sandia National Laboratories**

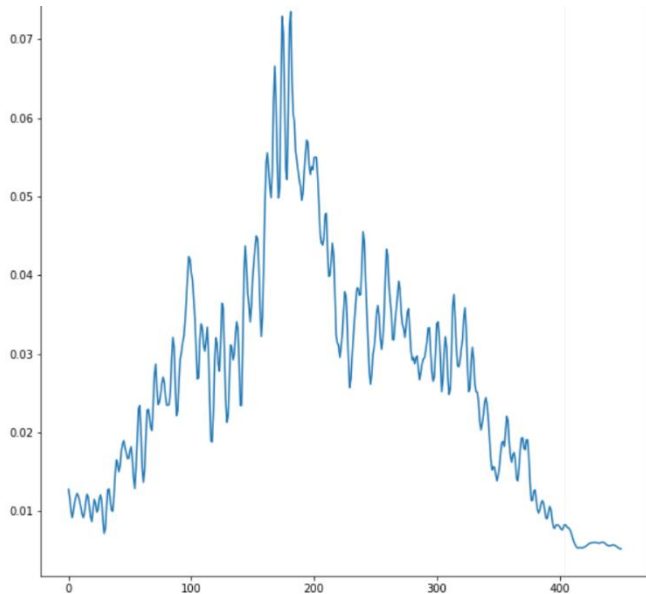**Exceptional service in the national interest**

# Classification of Optical Ports Using Machine Learning
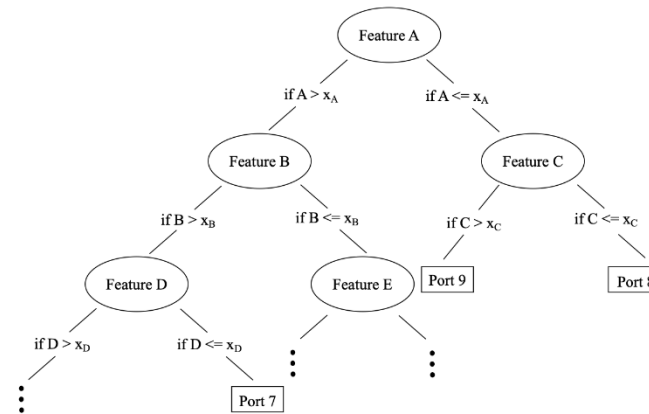
## Abigail Pribisova

MARTIANS End of Summer Symposium 2021

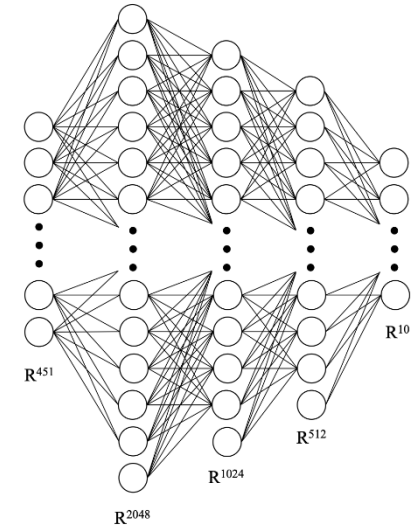U.S. DEPARTMENT OF **ENERGY**

NNSA

# Objective

- Predict the optical port from which a particular spectrum was produced
- Compare the effectiveness of decision trees versus fully-connected neural networks
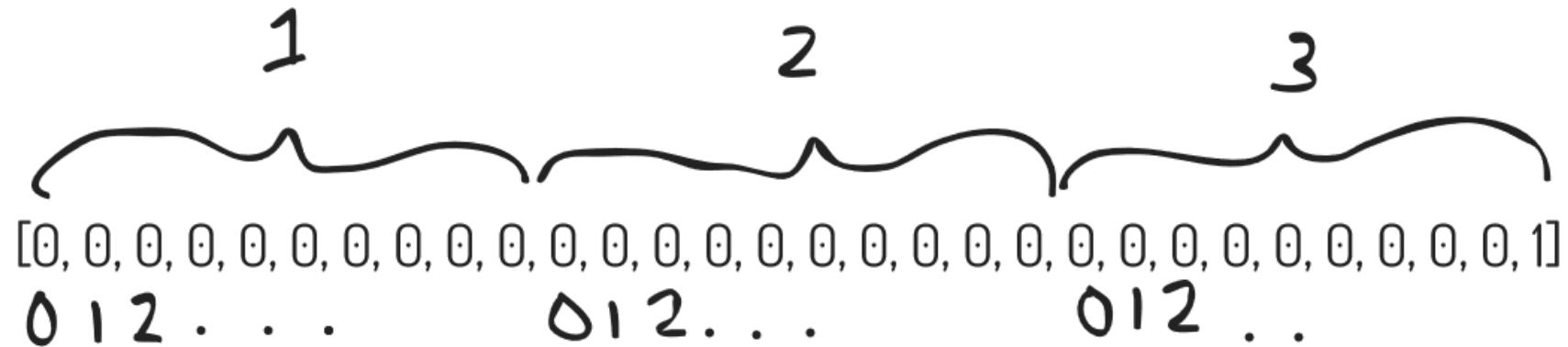
# Data Generation

- 750,000 spectrum-port pairs – 250,000 from each optical device
- Spectra = 451-component vectors
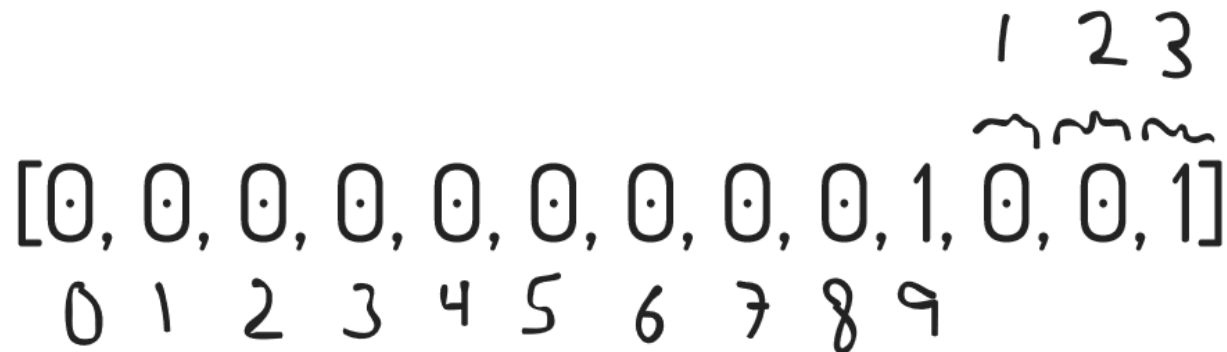- Ports (0-9) = label encoding ( [9] ) or one-hot vector encoding ( [0, 0, 0, 0, 0, 0, 0, 0, 0, 1] )

**70%**            **24%**        **6%**

# Data Generation cont.

- Multilabel classification -> classify both port number AND optical device
    1. 30-component vector



$[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1]$

    2. 13-component vector



$[0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1]$
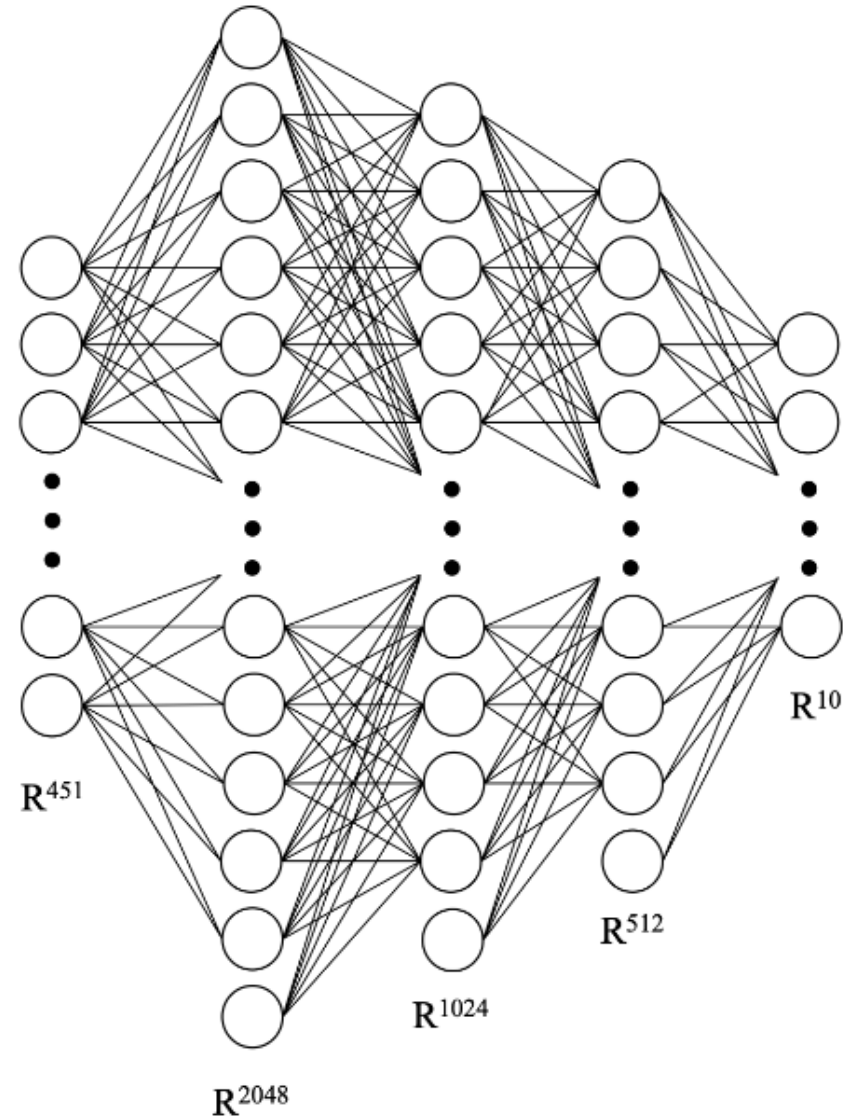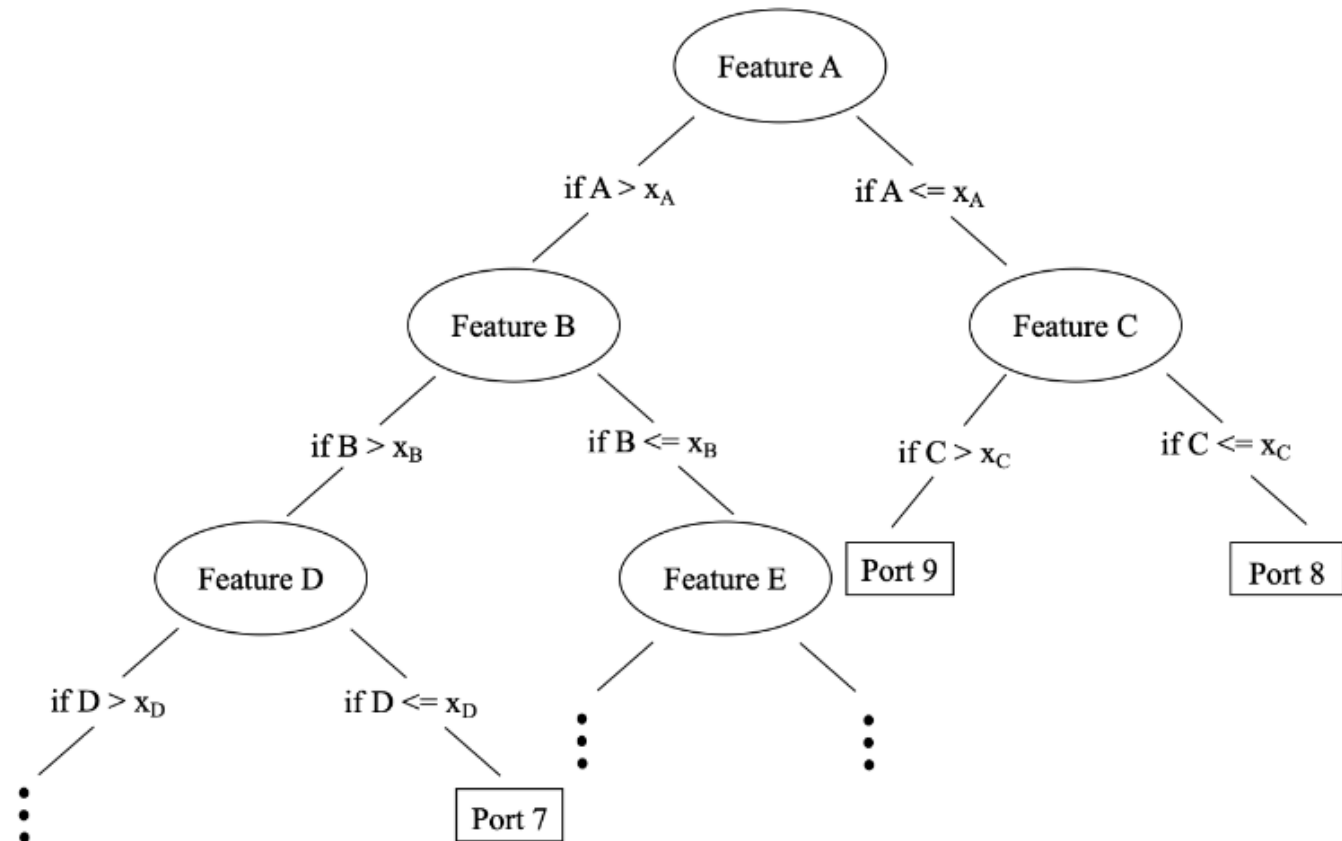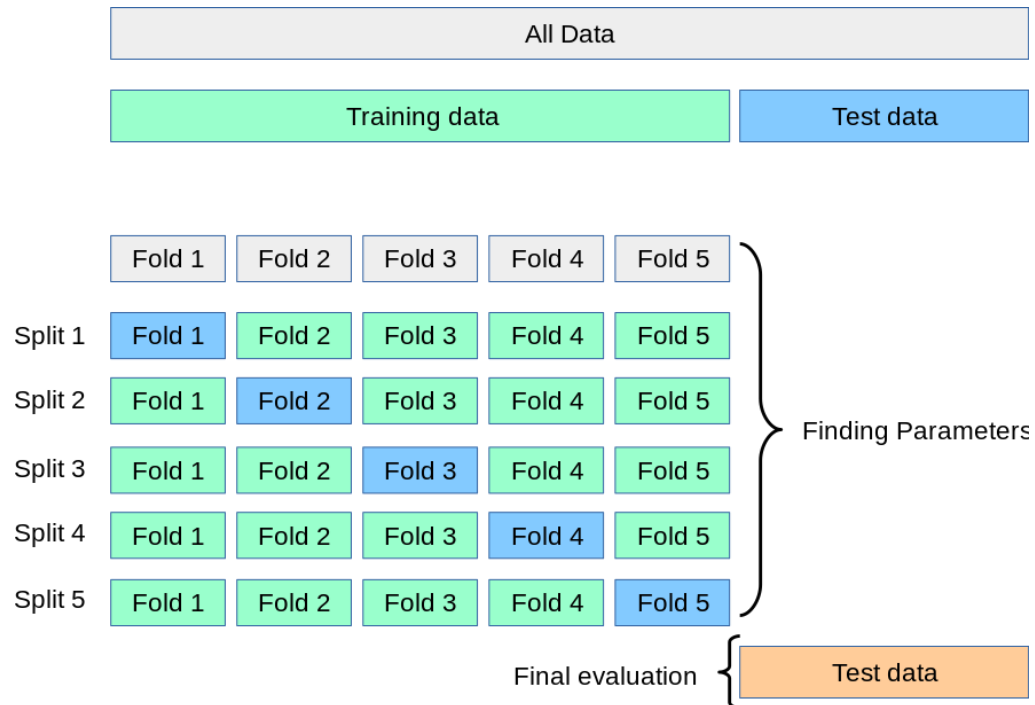
# Machine Learning Methods

- Inner activation function: **ReLU**
- Final activation function: **softmax**
- Optimizer: **Adam**
- Loss: **categorical crossentropy**
- Training time: **3 hours on 32 GB GPU**
- BatchNormalization layer
- EarlyStopping function



$R^{451}$   $R^{2048}$   $R^{1024}$   $R^{512}$   $R^{10}$
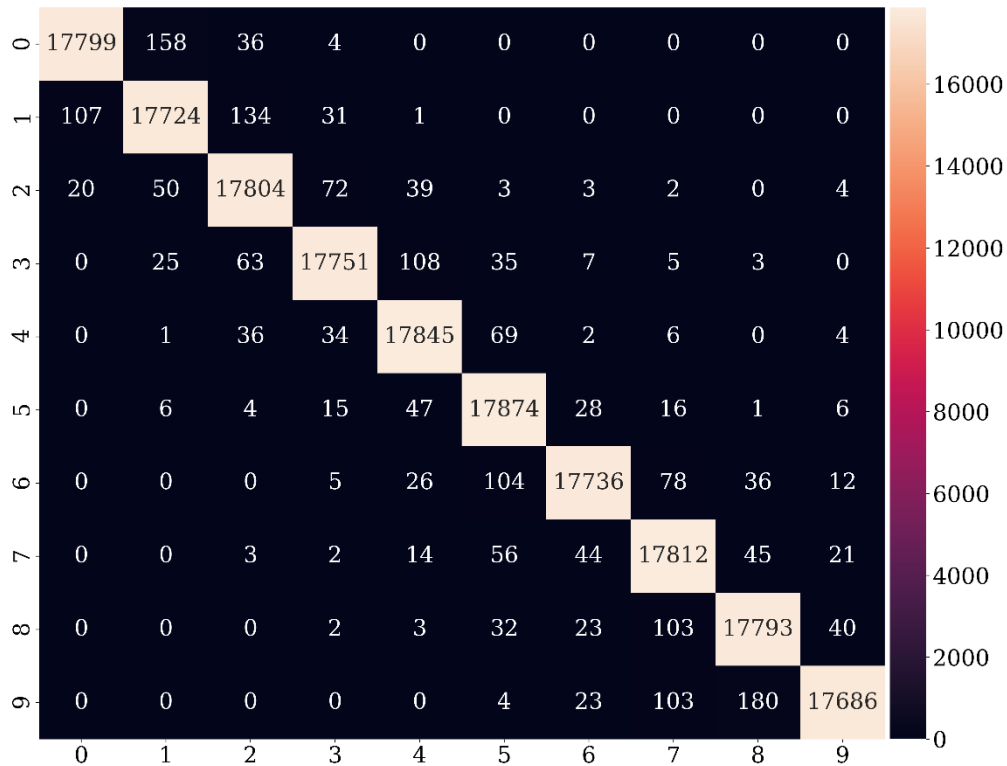
# Machine Learning Methods cont.

- scikit-learn DecisionTreeClassifier
- Criterion: **Gini impurity function**
- Minimum samples: **2**
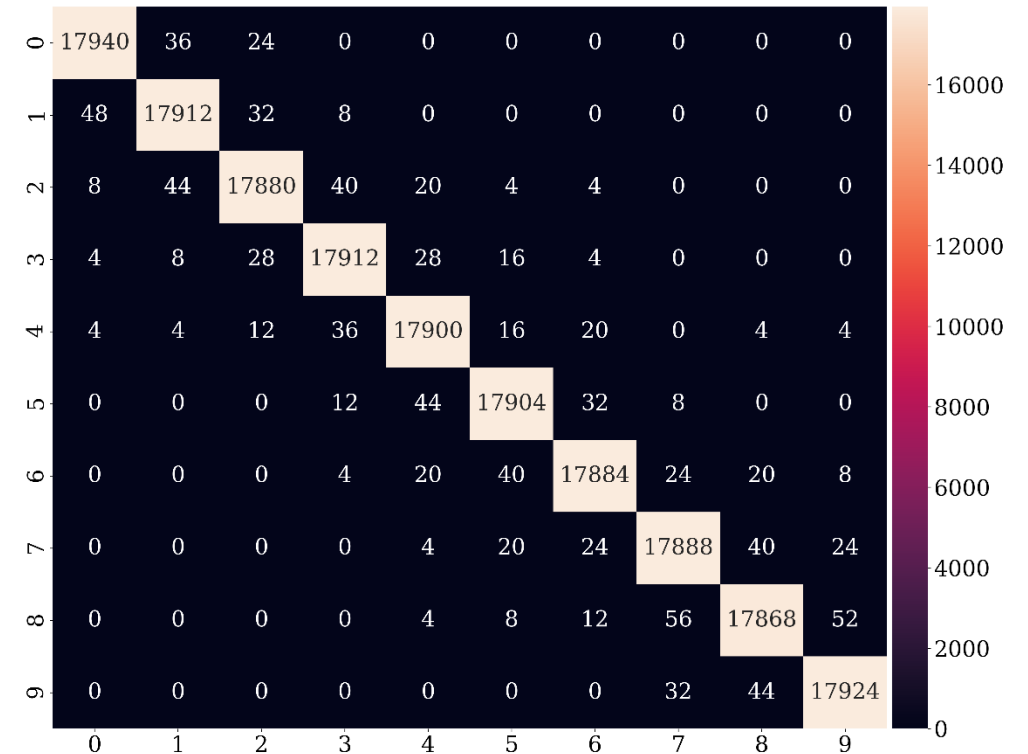- Training time: **1 hour on 1 CPU**

# Results

- Neural network: 98.07% (standard deviation: 0.727%)
- Decision tree: 99.43% (standard deviation: 0.0112%)
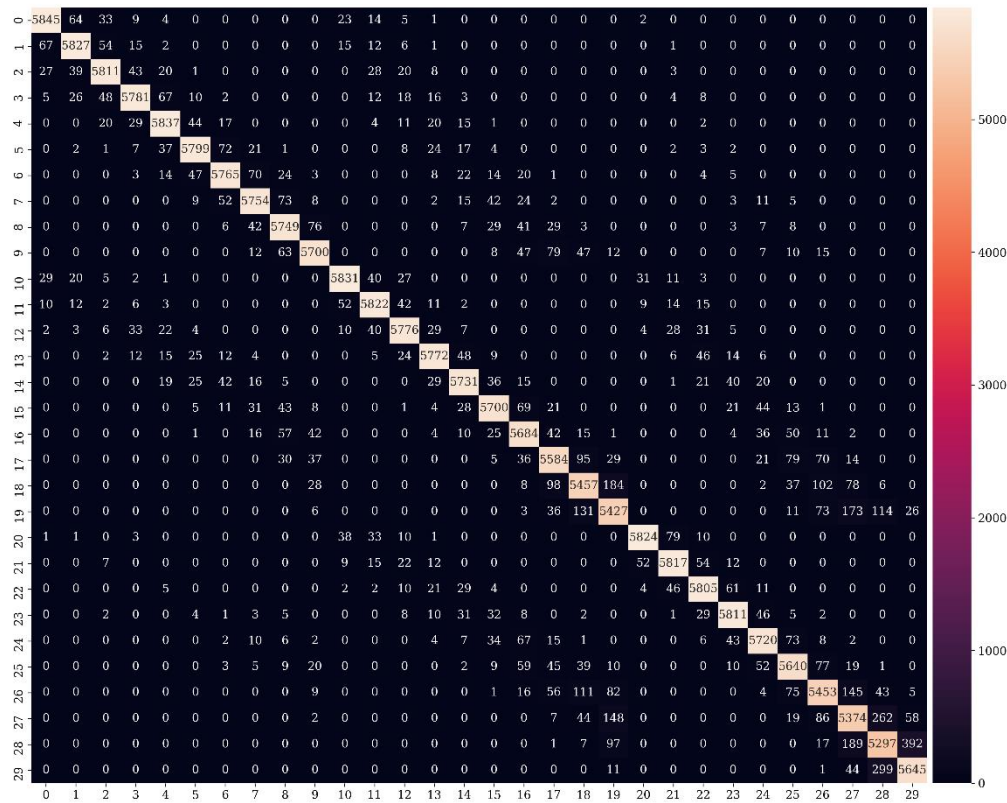


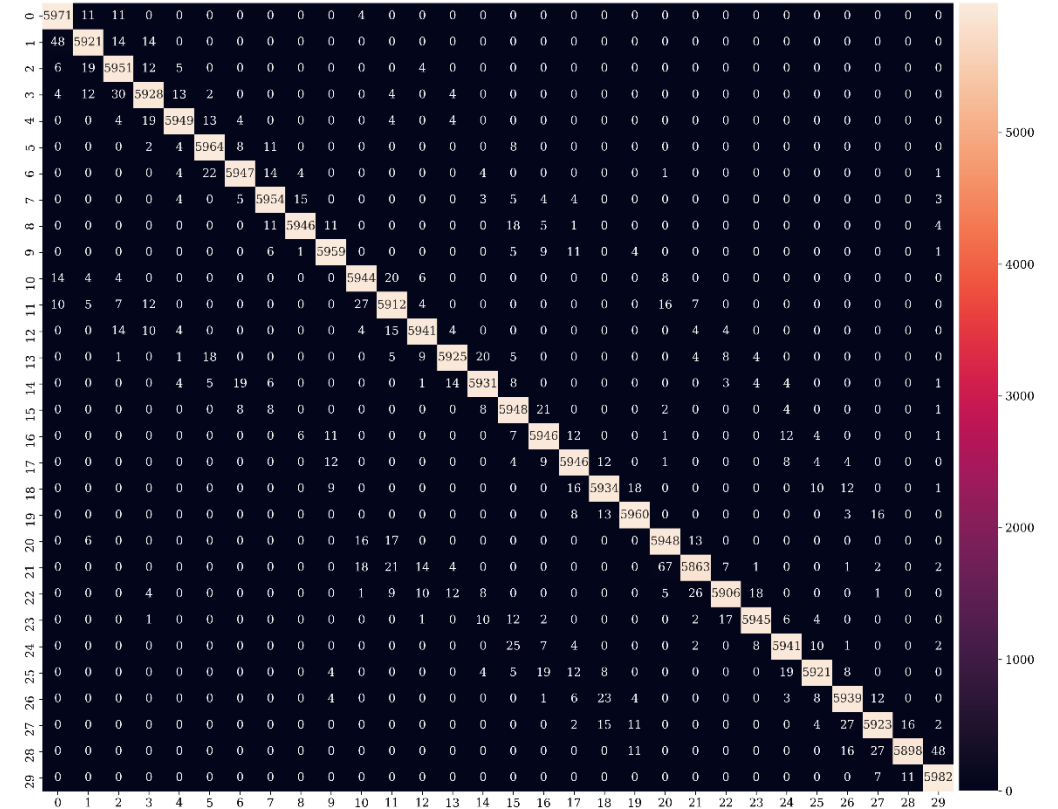Neural Network



Decision Tree

# Results cont.

- Multilabel neural network: 98.99%
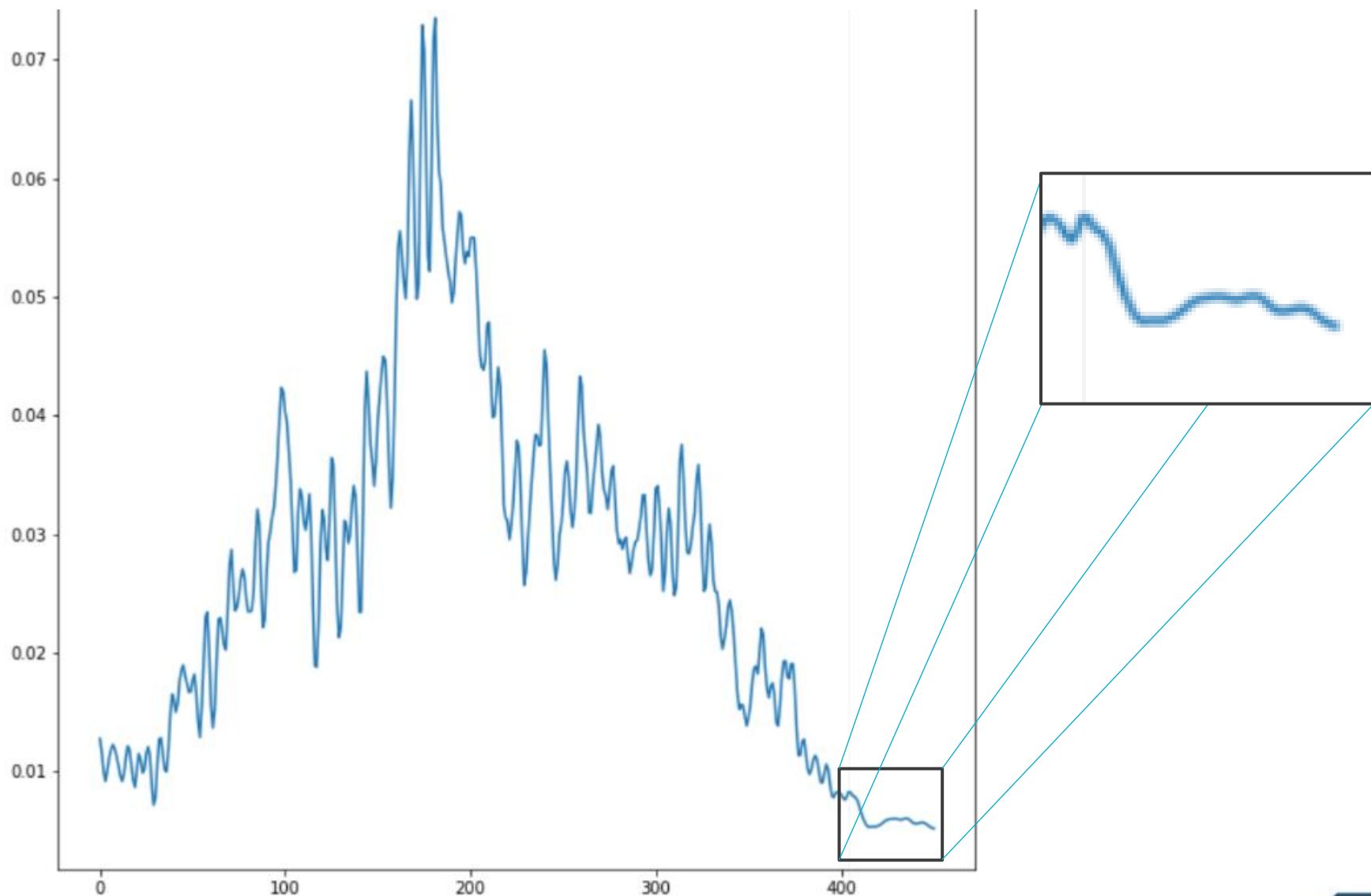- Multilabel decision tree: 95.02%



Neural Network



Decision Tree

# Results cont.

- Feature importance

# Results cont.

| Percentage of data / number of data samples used for training and testing (% / samples) | Fully-connected neural network accuracy (%) ~ standard deviation (%) over 10 runs | Decision tree classifier accuracy (%) ~ standard deviation (%) over 10 runs |
|---|---|---|
| 100 / 750,000 | 98.07 ~ 0.727 | 99.43 ~ 0.0112 |
| 90 / 675,000 | 96.81 ~ 1.054 | 98.88 ~ 0.0137 |
| 80 / 600,000 | 96.25 ~ 0.644 | 98.29 ~ 0.0157 |
| 70 / 525,000 | 93.36 ~ 1.512 | 96.54 ~ 0.0194 |
| 60 / 450,000 | 89.96 ~ 2.809 | 94.32 ~ 0.0395 |
| 50 / 375,000 | 82.60 ~ 2.077 | 86.98 ~ 0.0564 |
| 40 / 300,000 | 73.61 ~ 2.171 | 81.03 ~ 0.0659 |
| 30 / 225,000 | 61.73 ~ 0.508 | 62.56 ~ 0.0781 |
| 20 / 150,000 | 39.73 ~ 1.411 | 37.65 ~ 0.1391 |
| 10 / 75,000 | 35.82 ~ 6.332 | 44.04 ~ 0.1682 |

# Conclusion

- Decision trees
  - Less hyperparamter tuning
  - Higher accuracy with less data
  - Better intuition for performance
  - More flexible with output vector format

- Neural networks
  - Less sensitive to length of output vector representation