

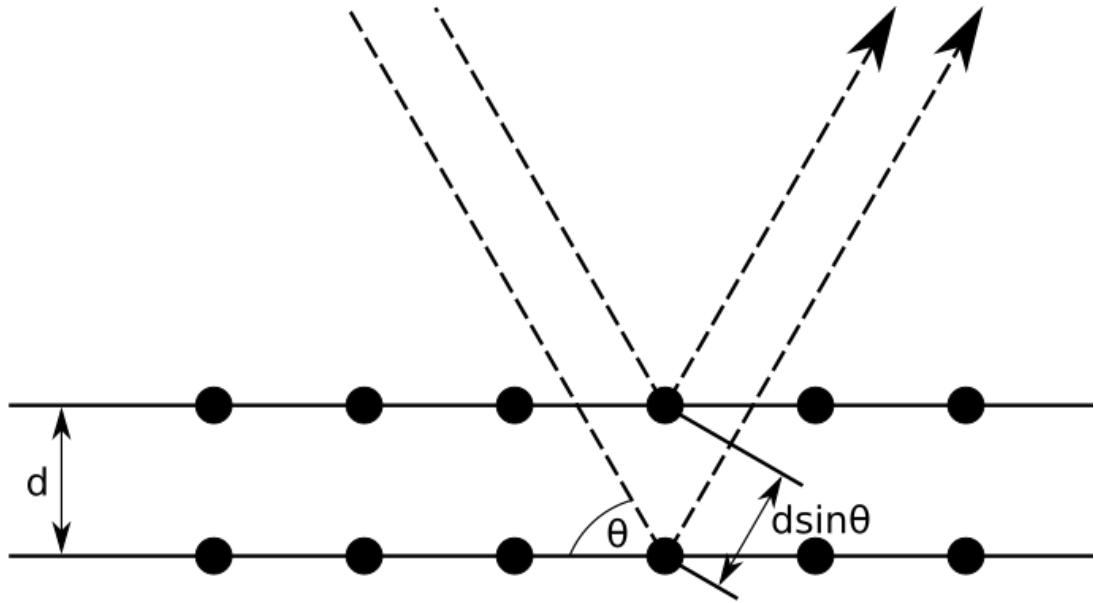
# Utilising Machine Learning on XRD Crystal Classification

侯奕安、吳政蔚、蔡杰達

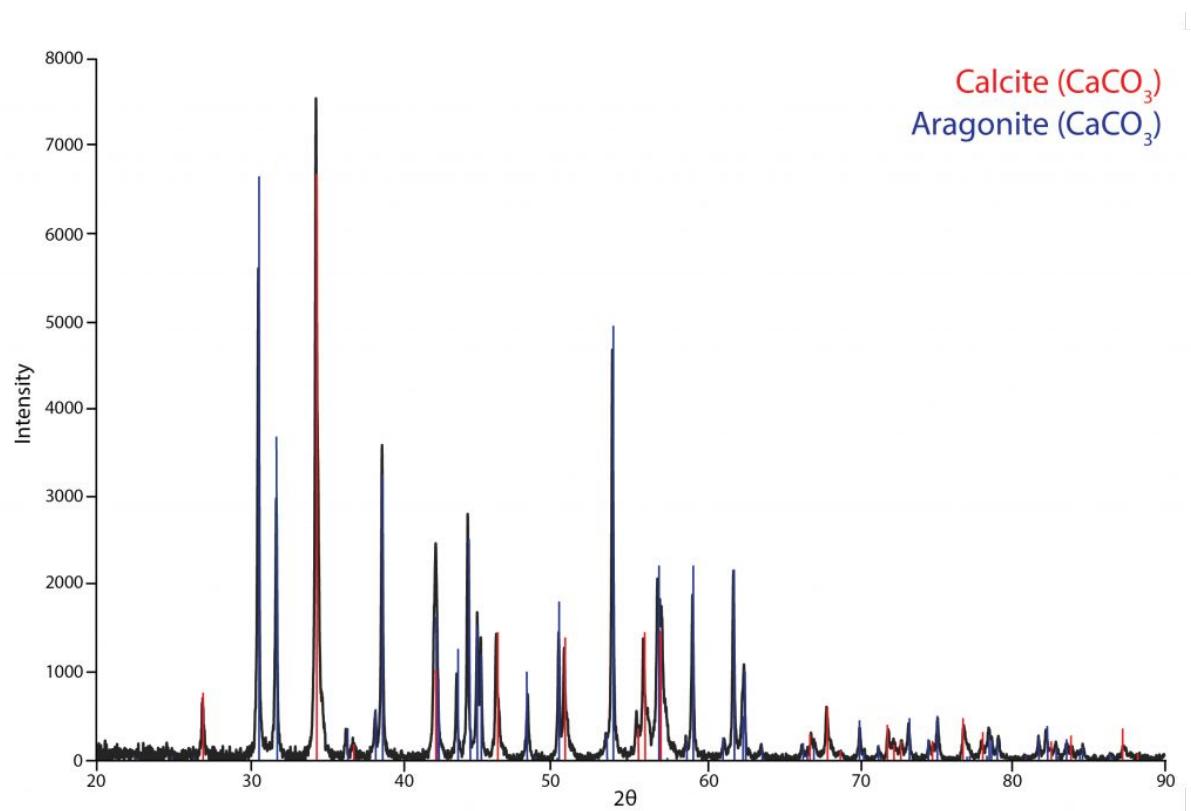
**What is XRD ?**

# Bragg's Law

$$2d\sin\theta = n\lambda$$



# XRD Spectrum



# Miller Indices

Consider a lattice plane with normal vector:

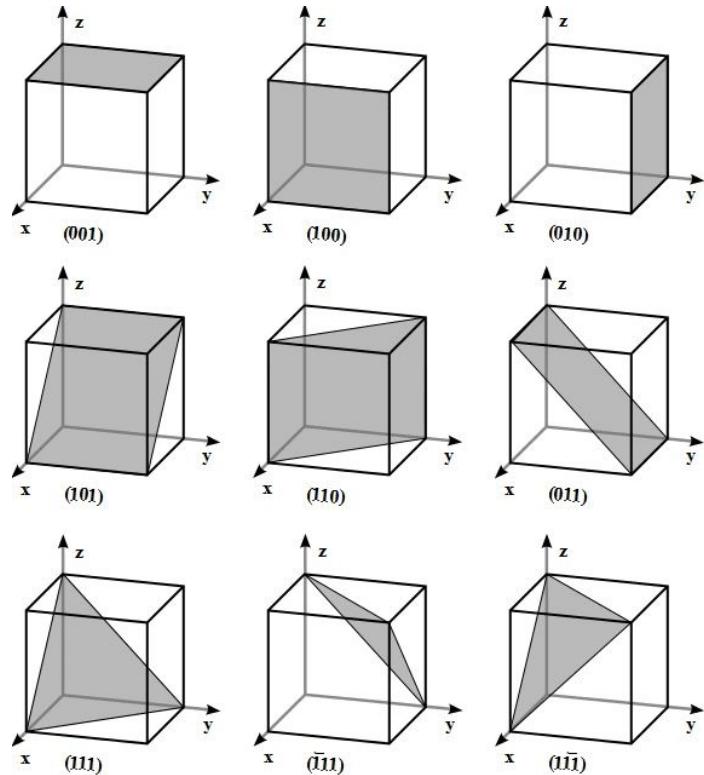
$$(1/h, 1/k, 1/l)$$

Then the Miller indices are denoted as:

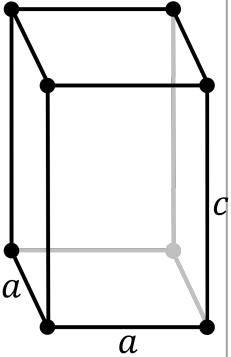
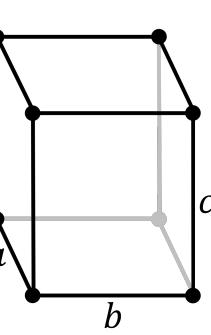
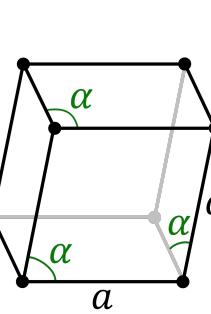
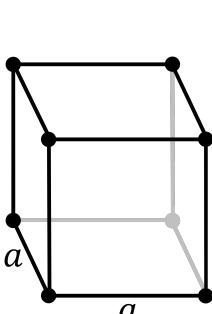
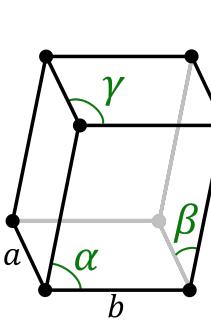
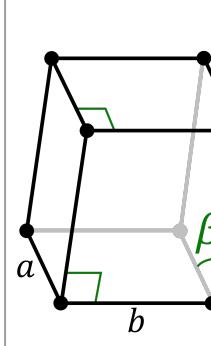
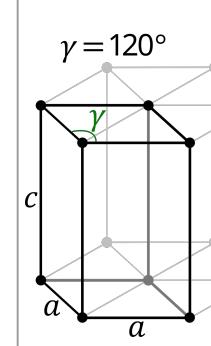
$$(h, k, l)$$

Classes are determined by relations between  $(h, k, l)$  and spacing between adjacent lattice plane.

Eg. cubic, hexagonal, etc

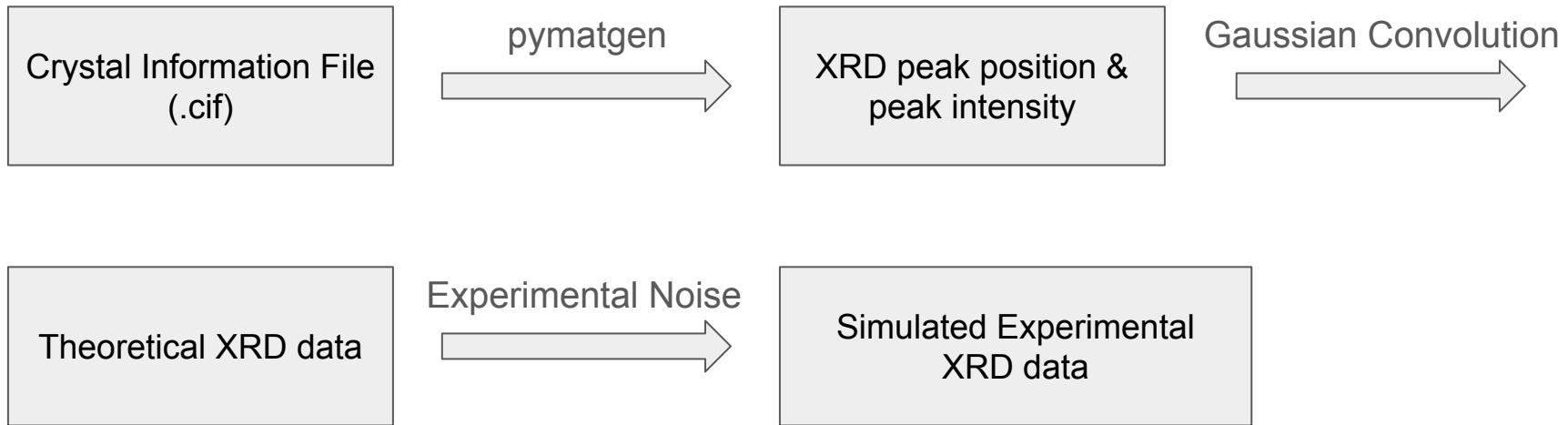


# Lattices

Tetragonal(0)	orthorhombic (1)	trigonal(2)	cubic(3)	triclinic(4)	monoclinic(5)	hexagonal(6)
						

# Data Preparation

# Data preparation



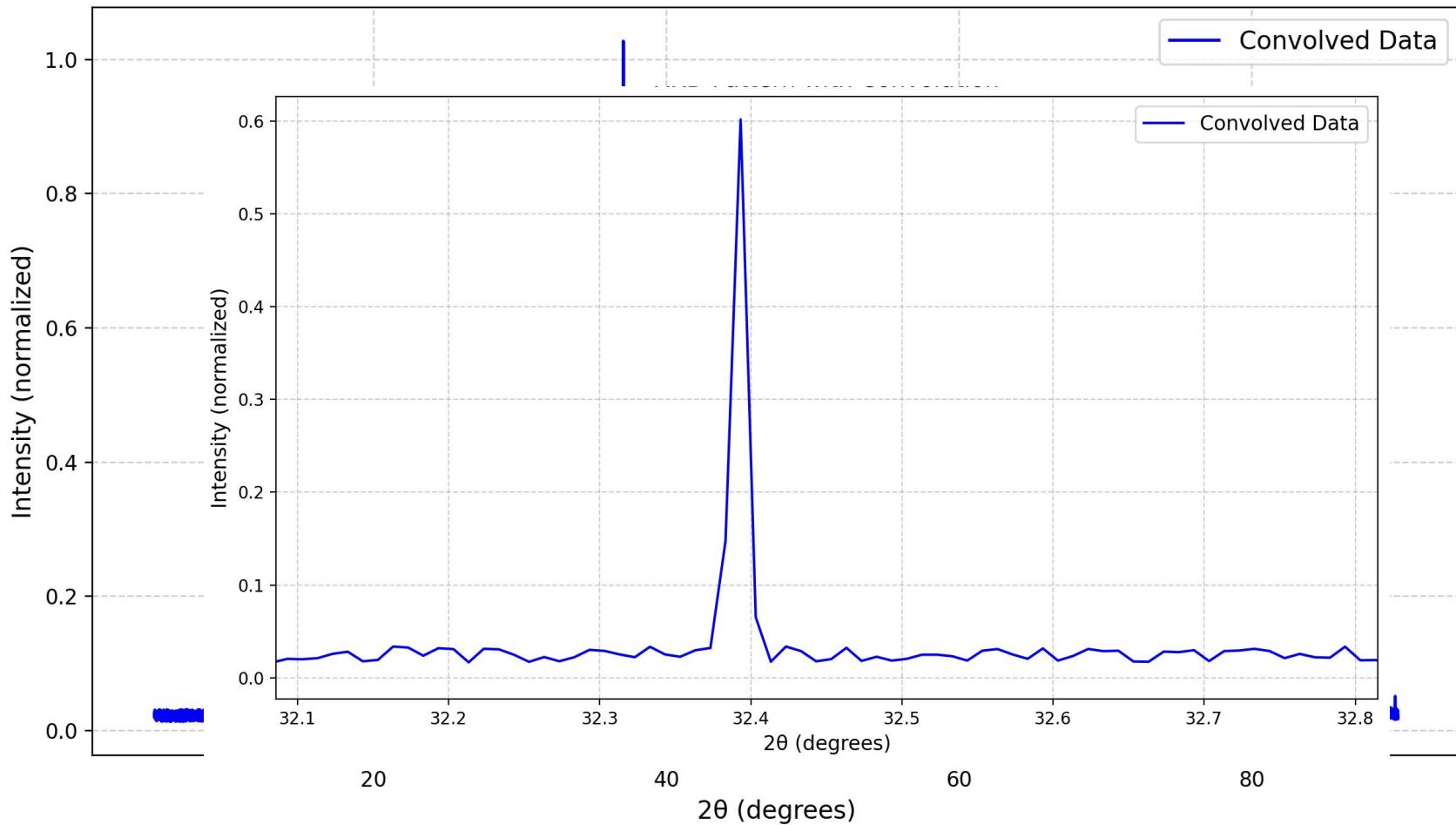
# Data preparation (CIF)

```
_chemical_formula_sum          'As Gd Te'  
_chemical_formula_weight       359.77  
_chemical_name_systematic  
;  
Gadolinium Asernic Telluride  
;  
_space_group_IT_number         62  
_symmetry_cell_setting        orthorhombic  
_symmetry_space_group_name_Hall '-P 2ac 2n'  
_symmetry_space_group_name_H-M  'P n m a'  
_atom_sites_solution_hydrogens geom  
_atom_sites_solution_primary   direct  
_atom_sites_solution_secondary difmap  
_audit_creation_method        SHELXL-97  
_cell_angle_alpha              90.00  
_cell_angle_beta               90.00  
_cell_angle_gamma              90.00  
_cell_formula_units_Z          4  
_cell_length_a                 7.5611(15)  
_cell_length_b                 4.0510(8)  
_cell_length_c                 9.920(2)  
_cell_measurement_temperature  153(2)  
_cell_volume                   303.85(10)
```

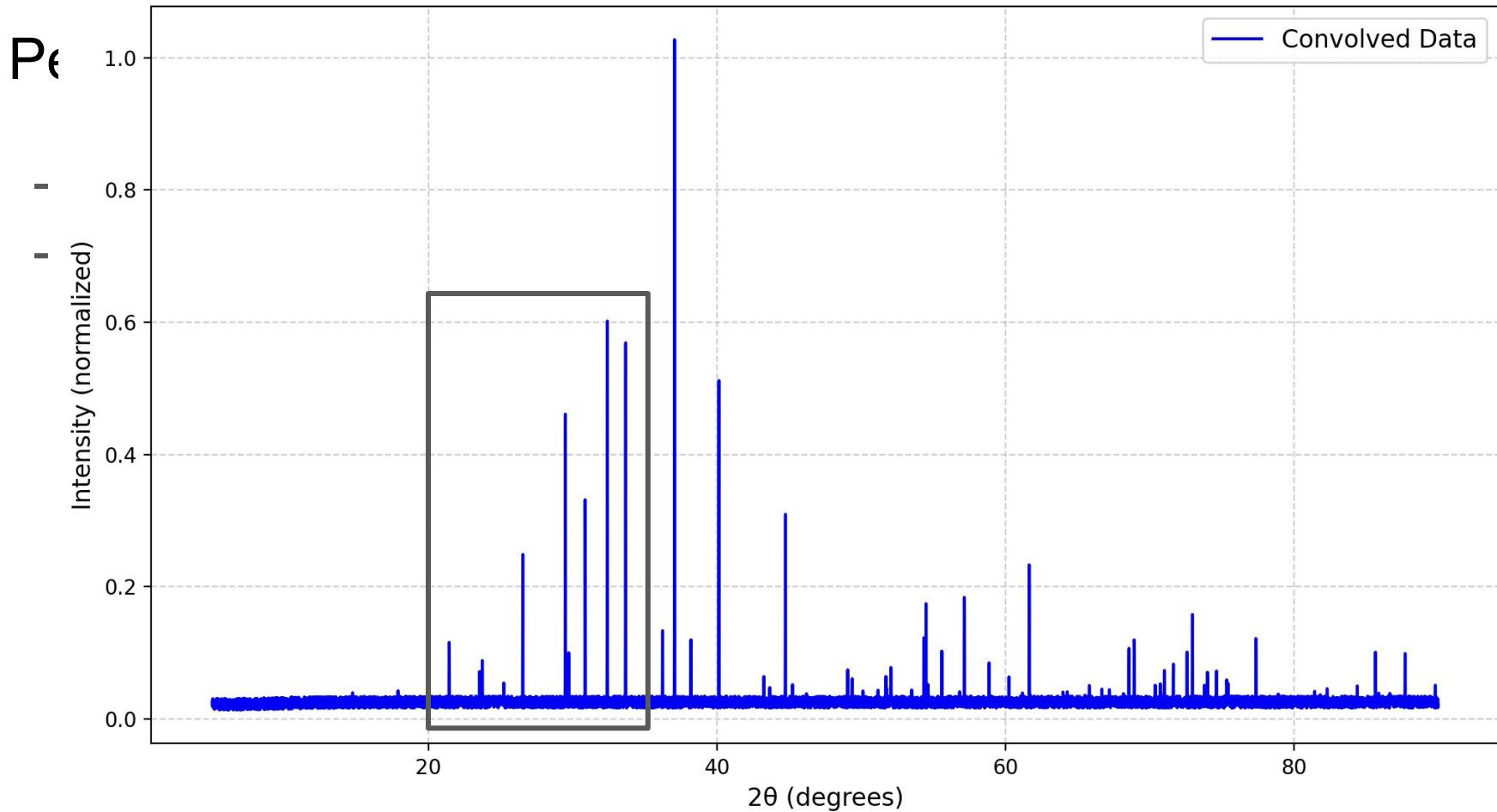
- We remove file with large crystal size (50 Å) and organic substance
- There are 19,525 data in total

# XRD Pattern with Convolution

D:



# XRD Pattern with Convolution

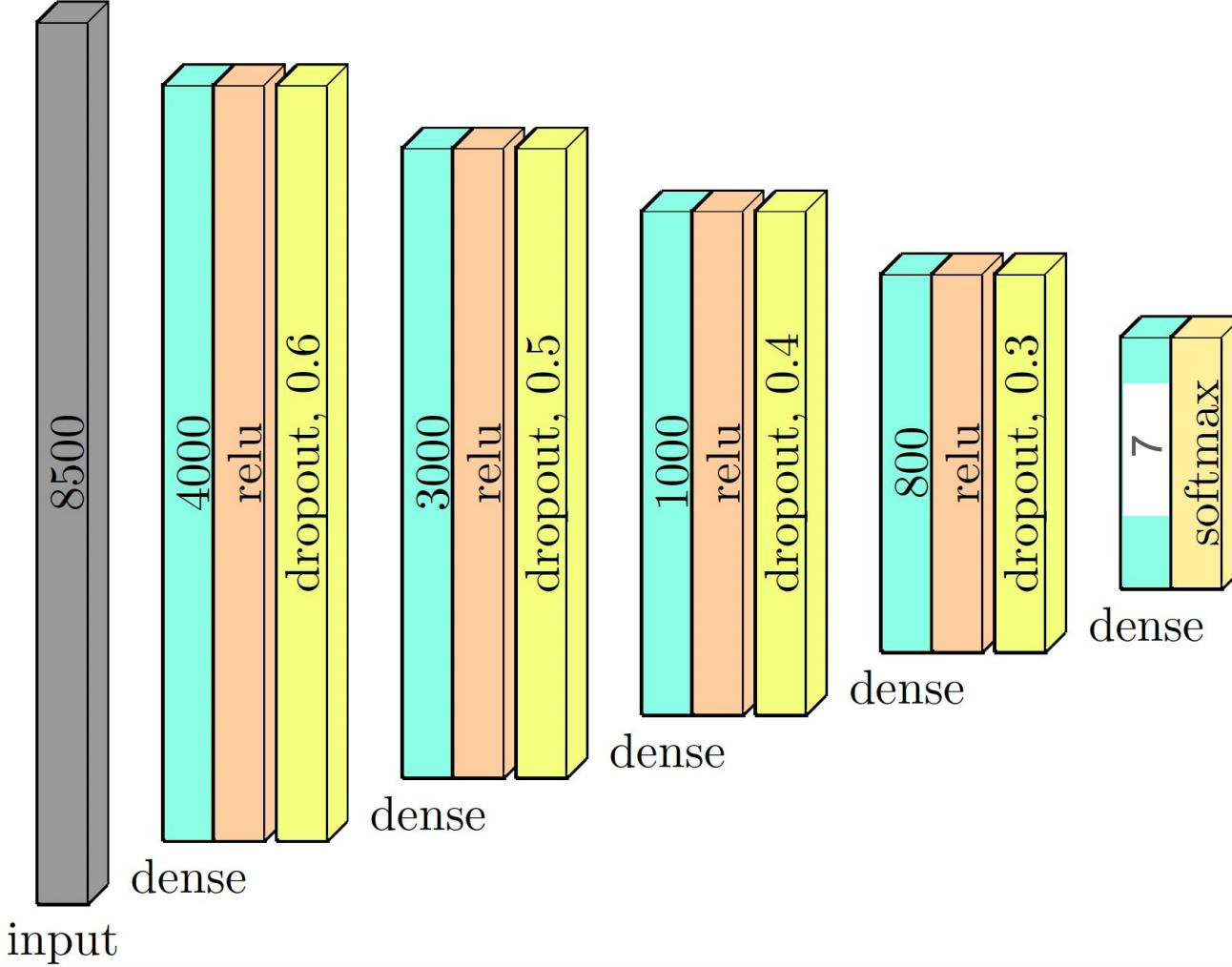


## Data train/test split

```
# Split data based on split_info
split_info = {
    "tetragonal": (3000, 500, 100),
    "orthorhombic": (3500, 500, 300),
    "trigonal": (2000, 500, 100),
    "cubic": (2000, 500, 100),
    "triclinic": (1800, 500, 10),
    "monoclinic": (3200, 500, 300),
    "hexagonal": (100, 10, 5)
}
```

# 1. Dense

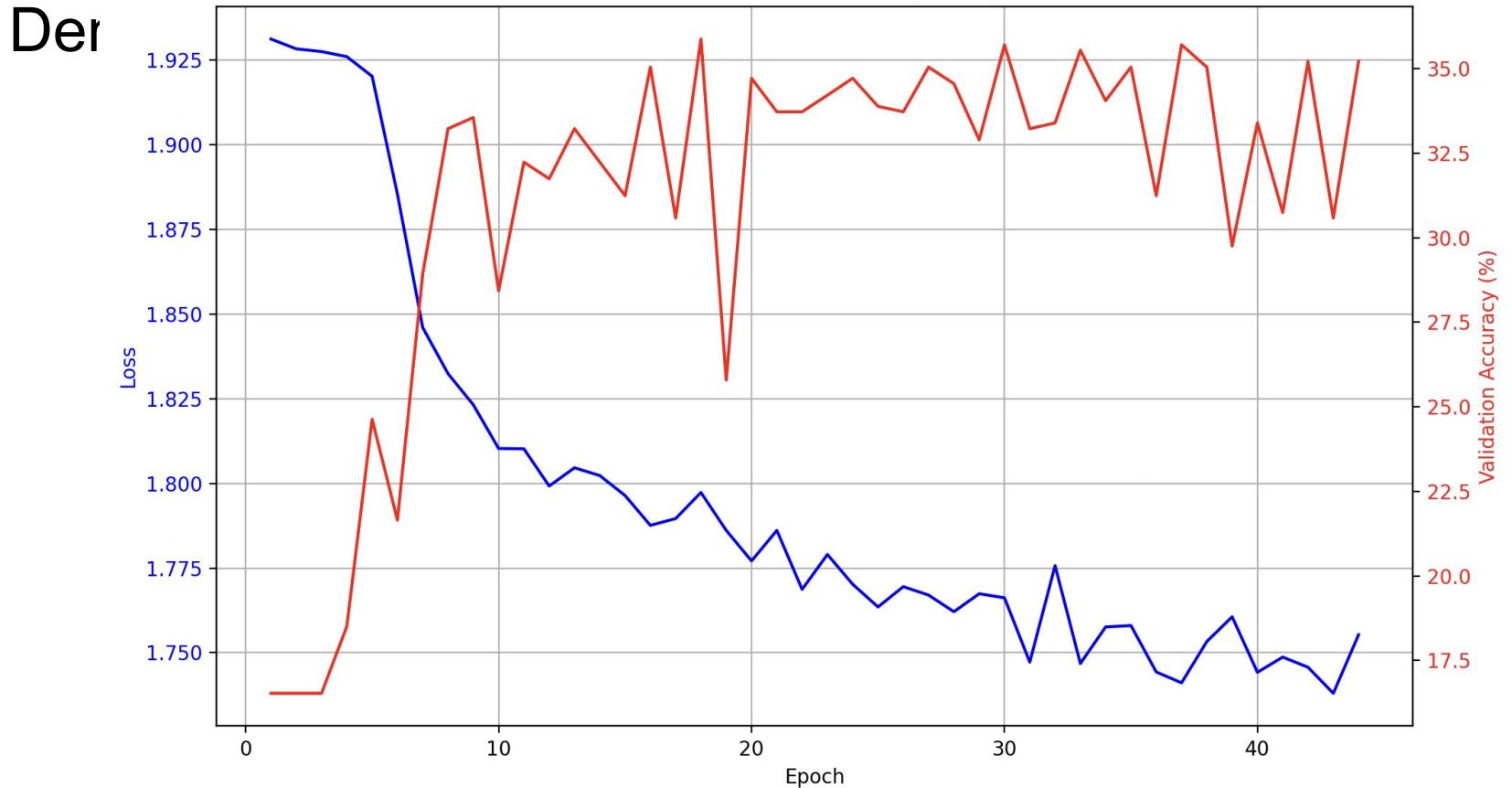
Dense



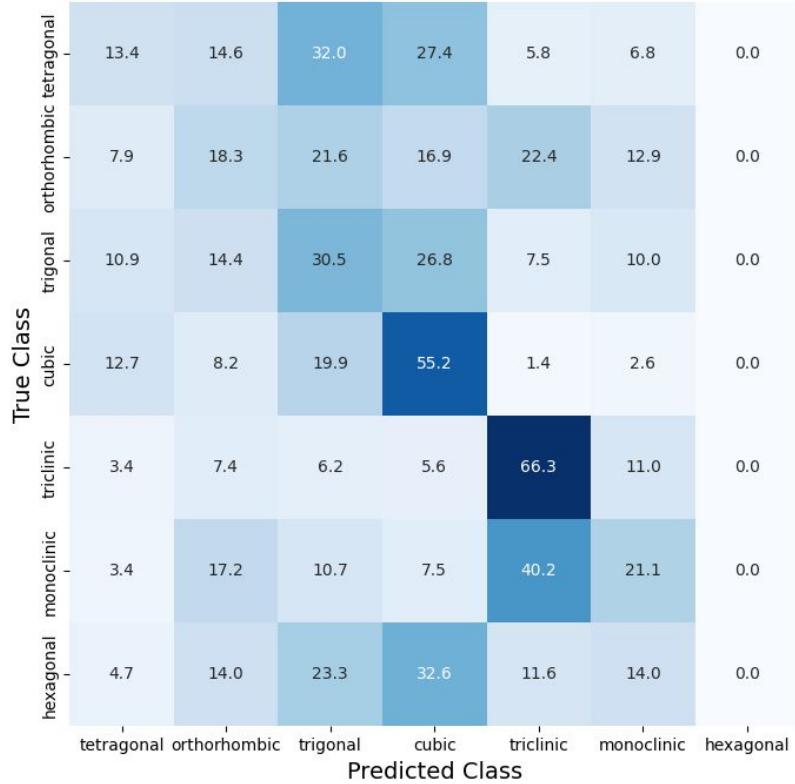
# Dense Model Accuracy

**All data** 35.25 %    **Peak data** 56.4 %

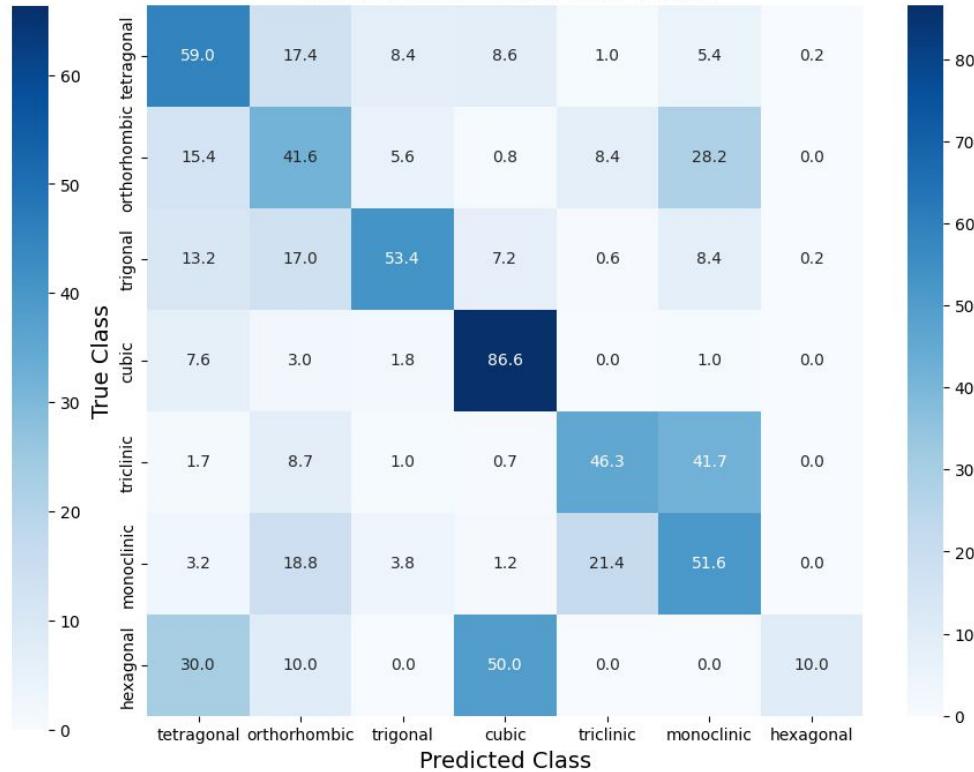
Loss and Validation Accuracy vs. Epoch



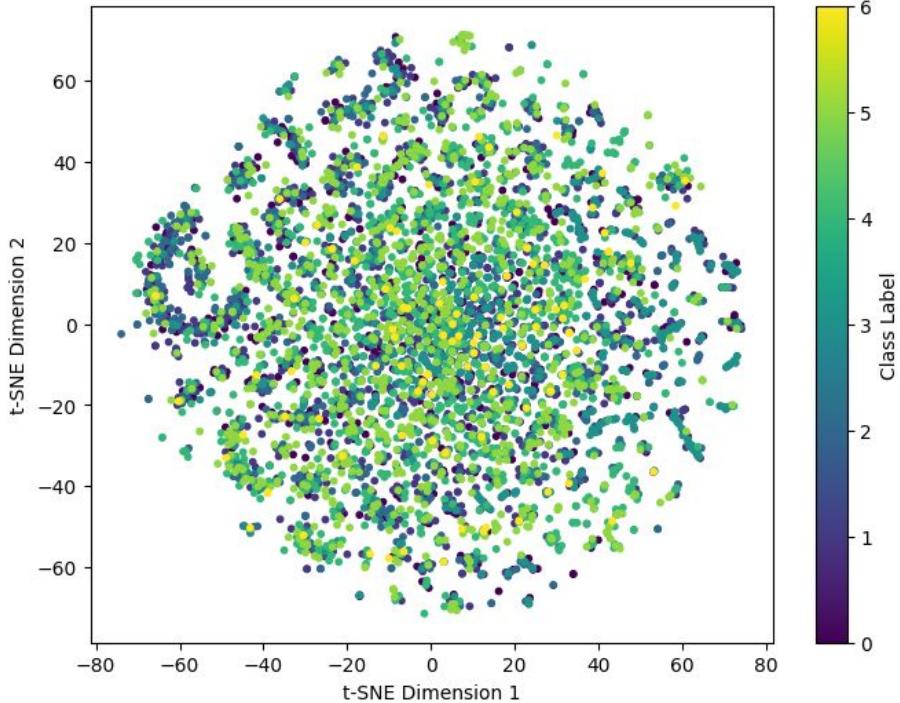
Confusion Matrix (Percentage Accuracy)



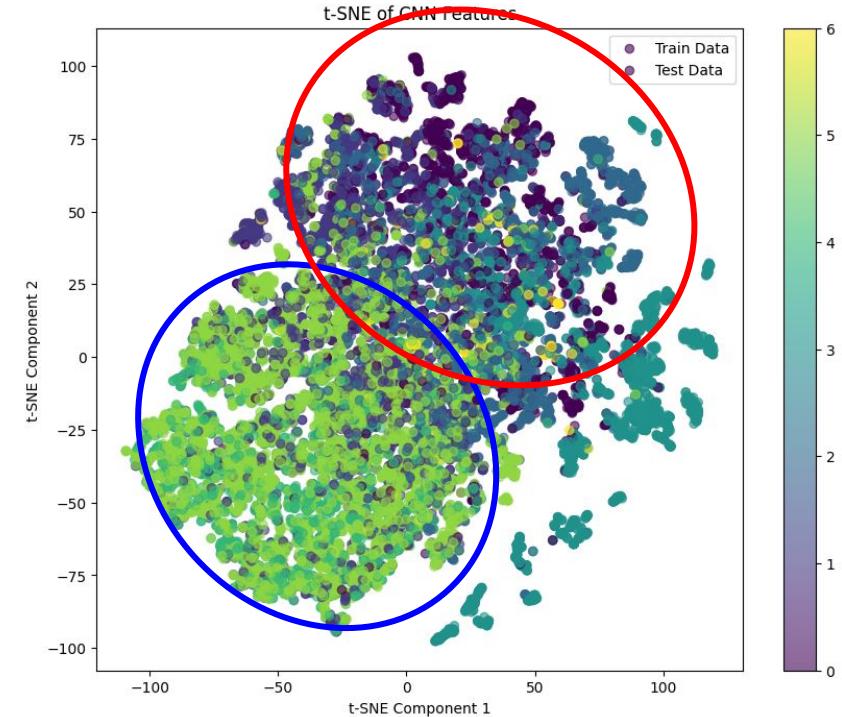
Confusion Matrix for CNN Model



t-SNE Visualization of XRD CNN Features

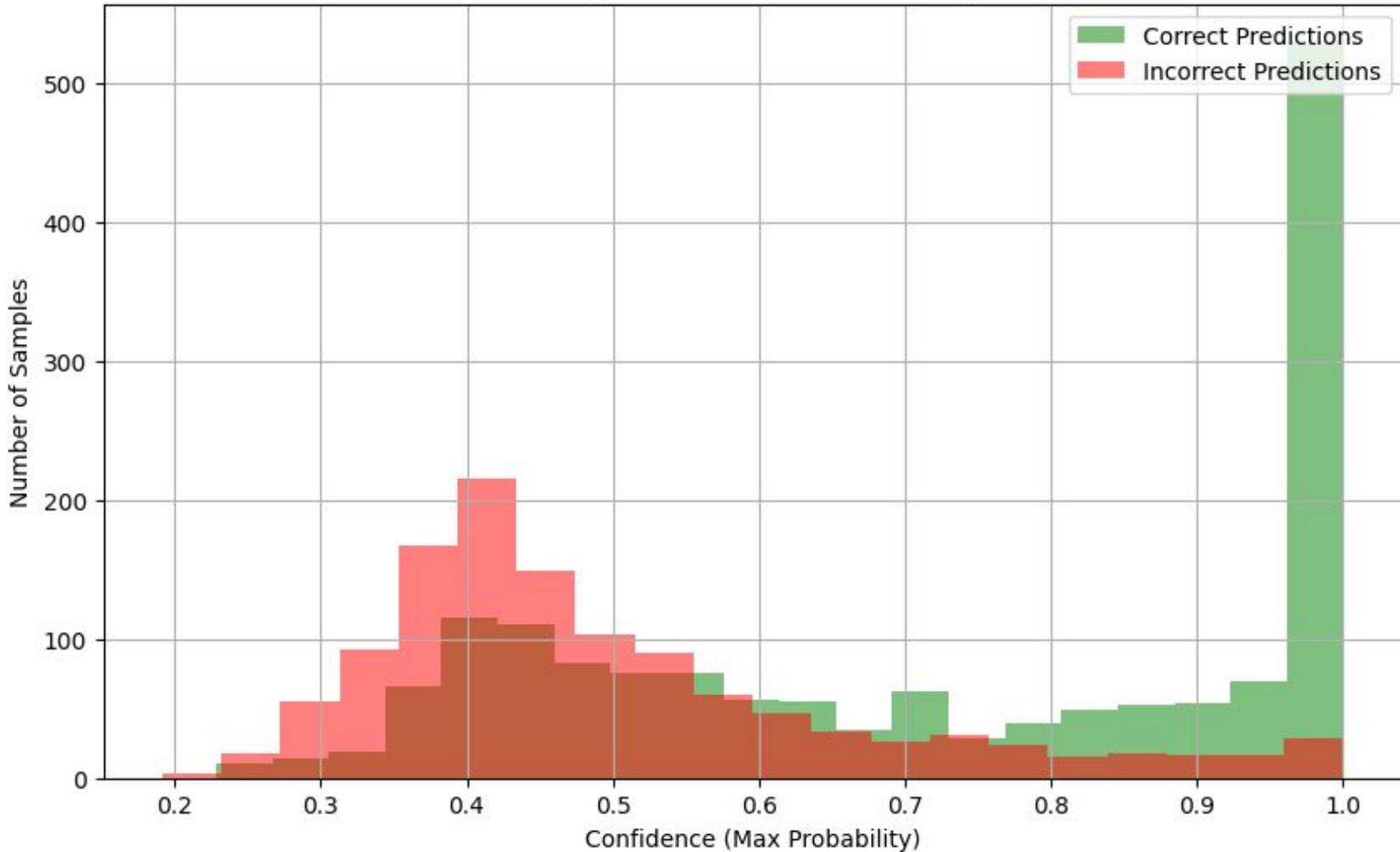


t-SNE of CNN Features



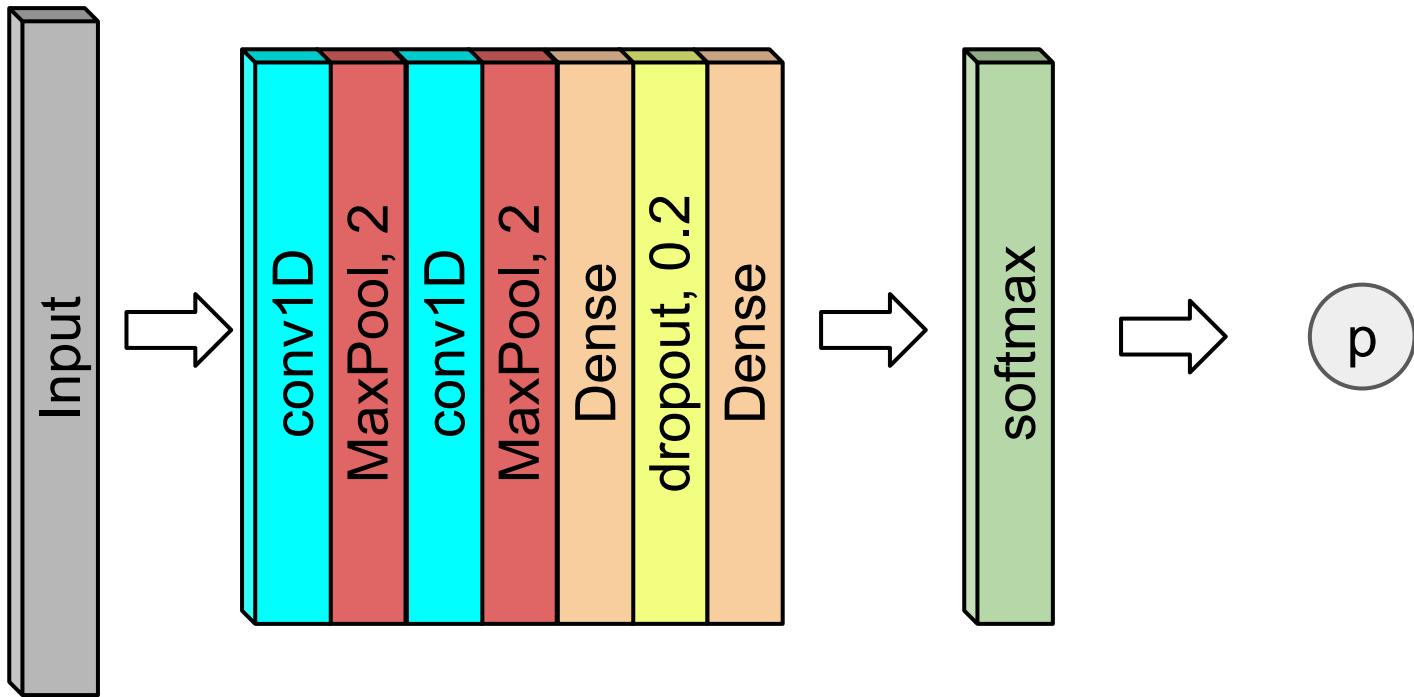
### Histogram of Prediction Certainty

D



## 2. CNN + Dense

# CNN-Dense: Architecture

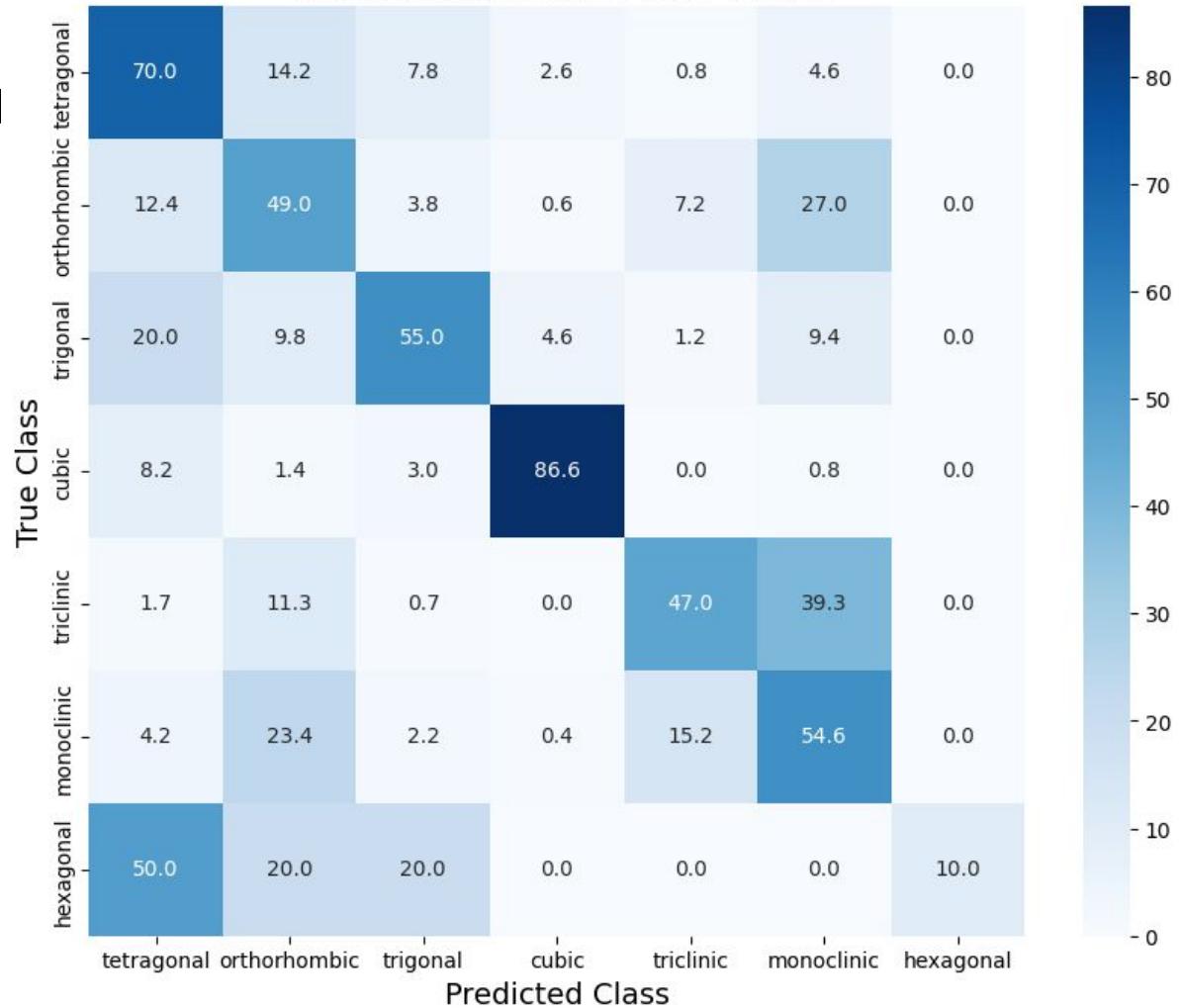


# CNN + Dense Model Accuracy

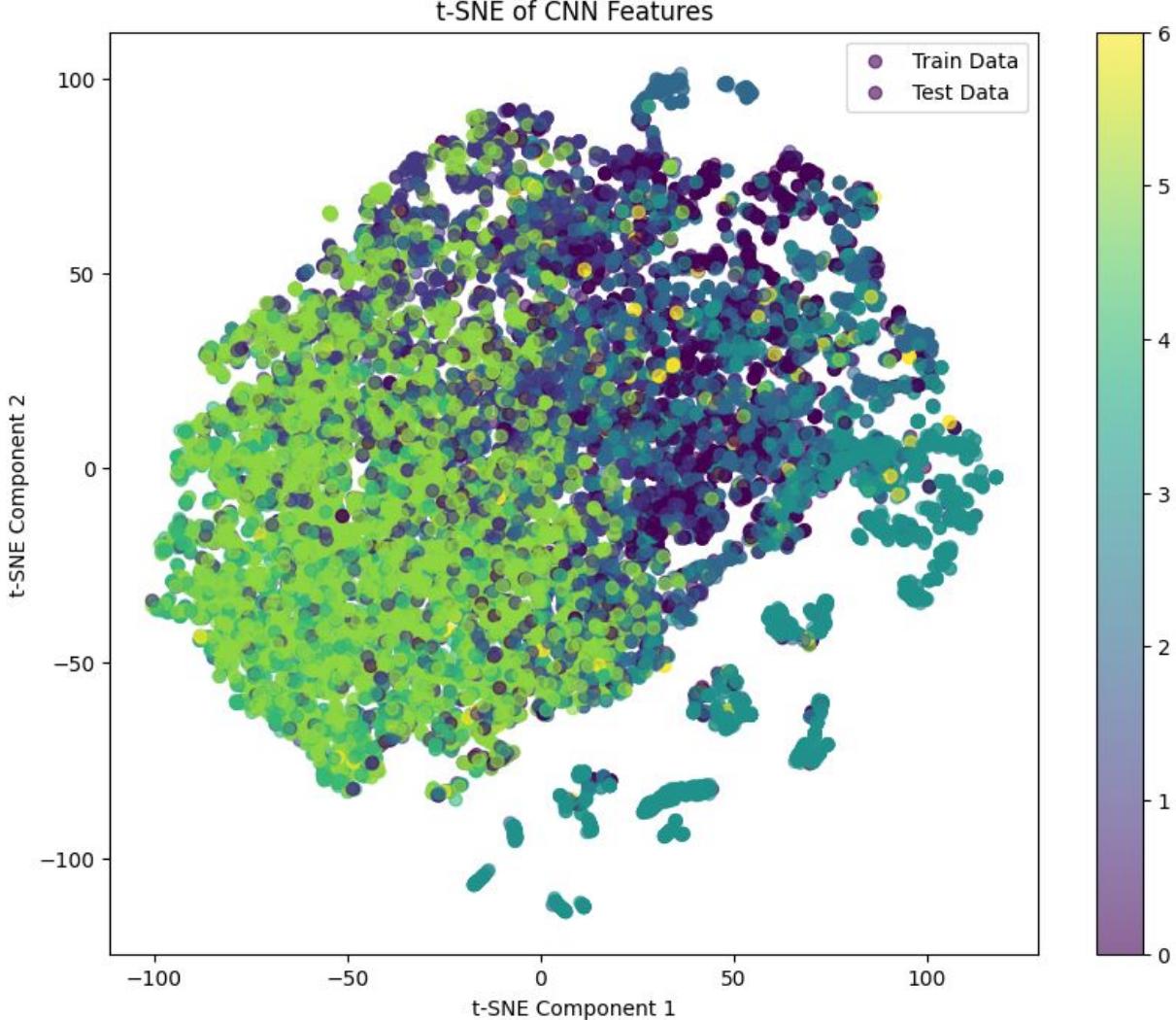
**Peak data 60.37 %**

Confusion Matrix for CNN Model

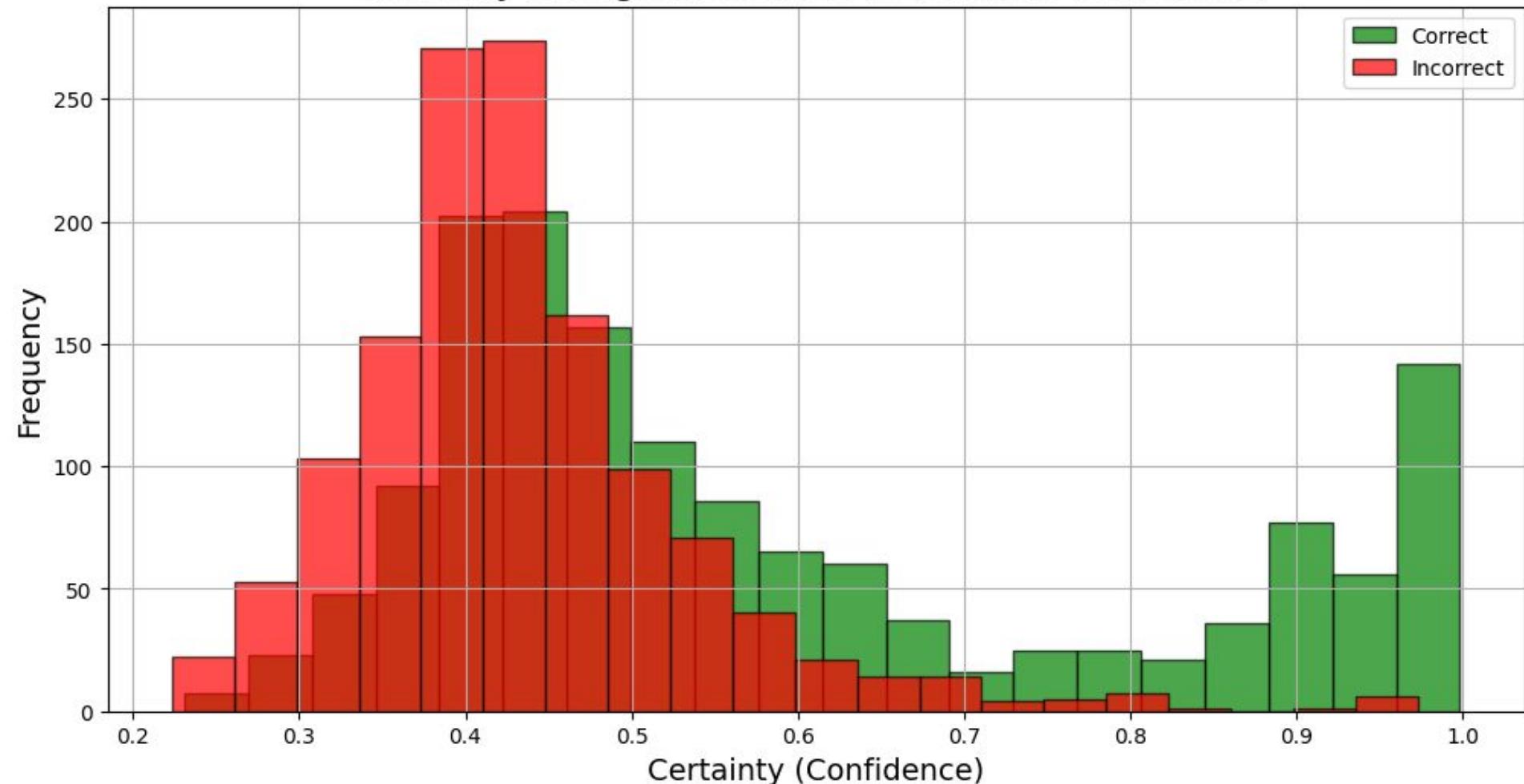
CNN+De



# CNN+De

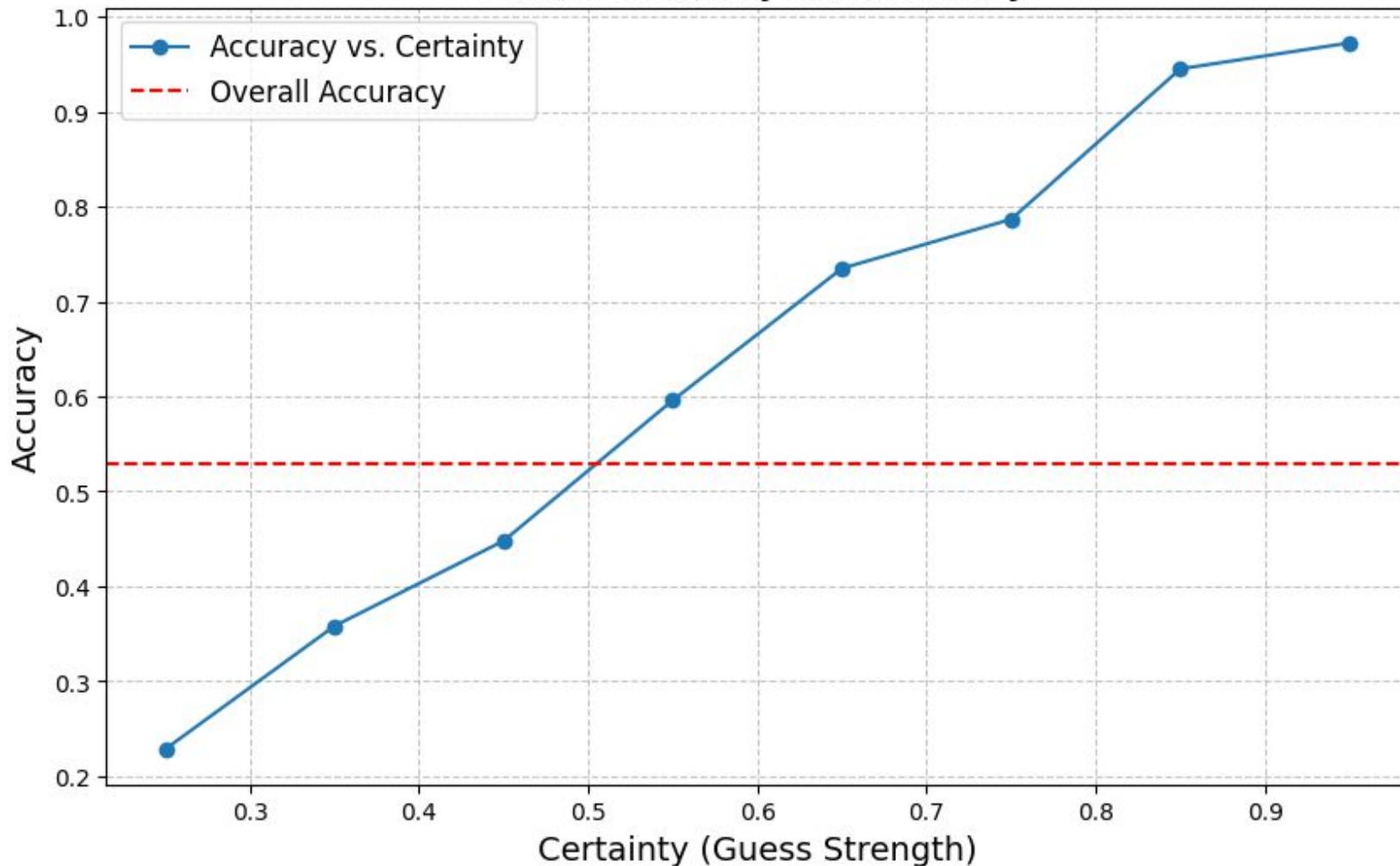


## Certainty Histogram (Correct vs Incorrect Predictions)

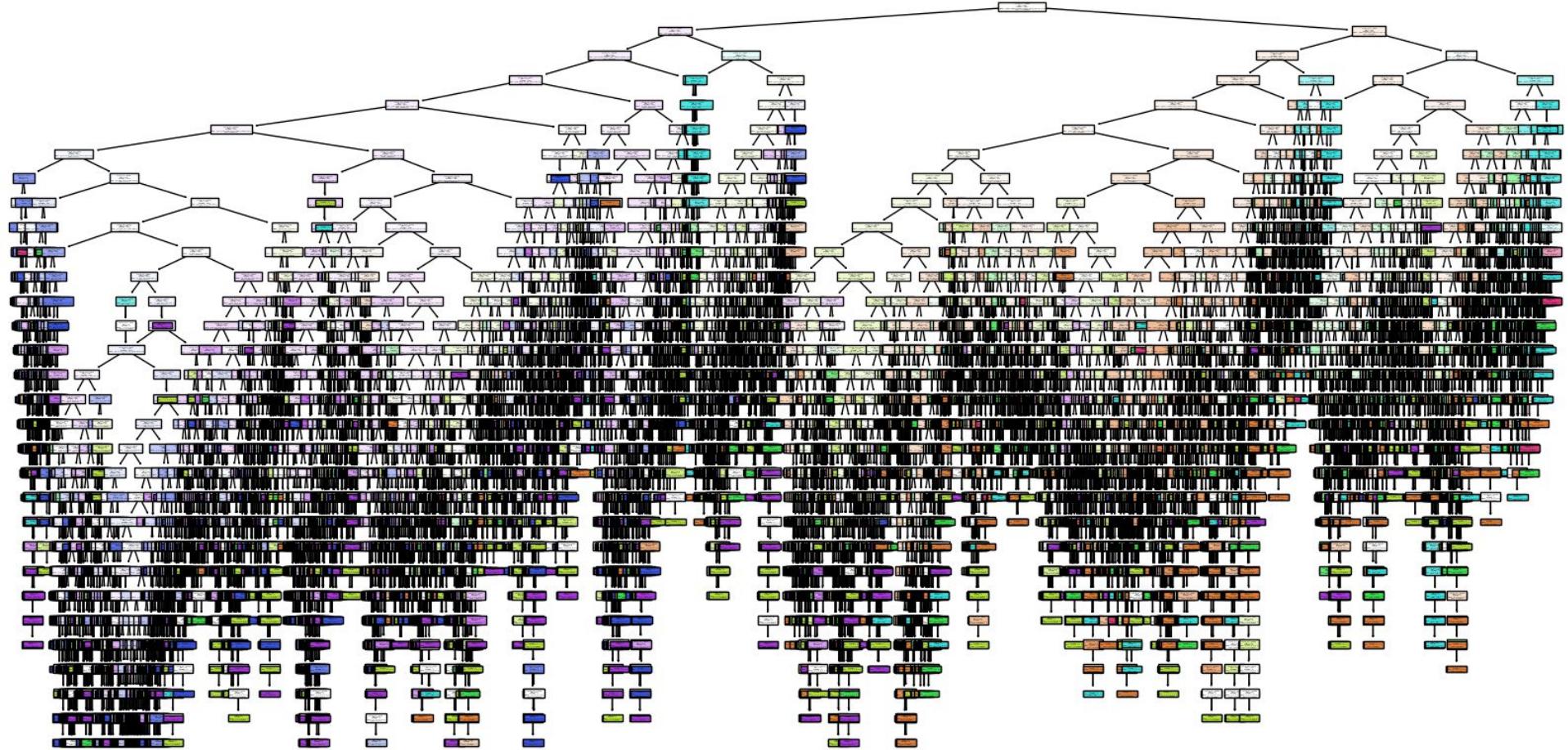


## Model Certainty vs. Accuracy

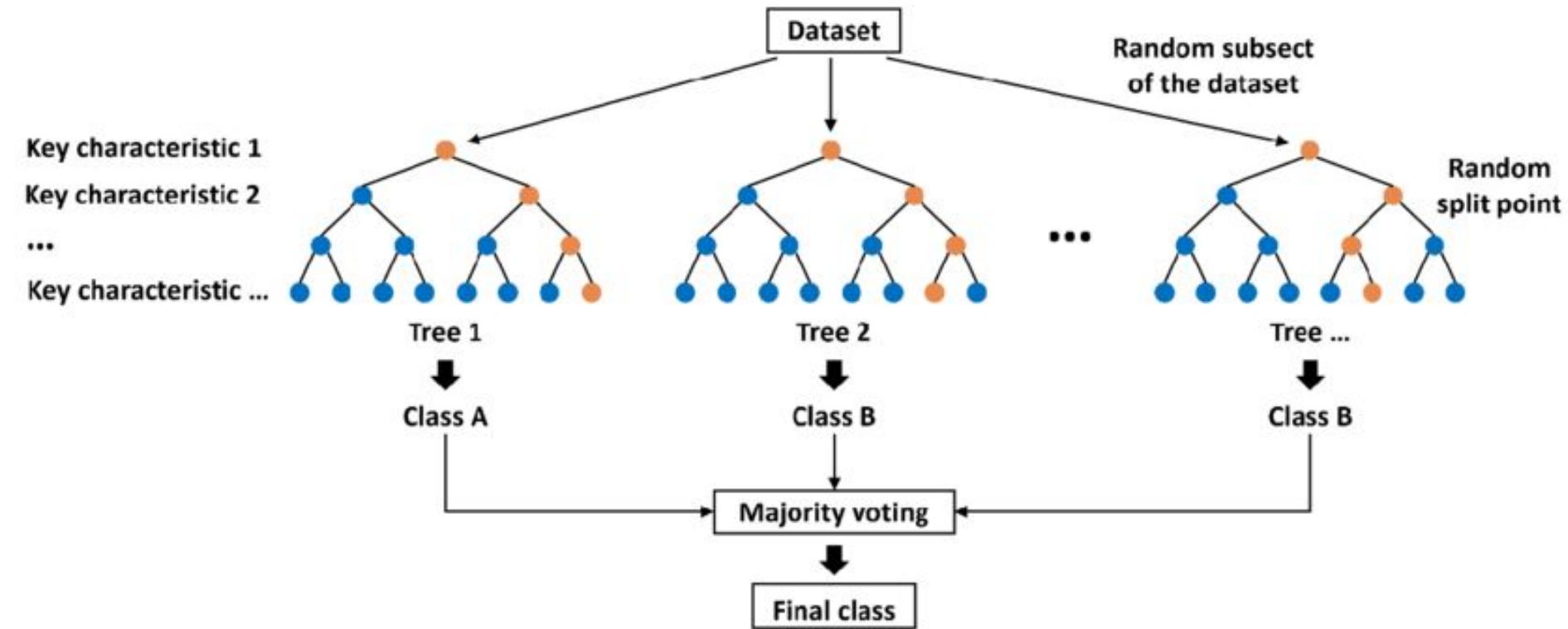
Ex



### 3. Extreme Randomized Trees



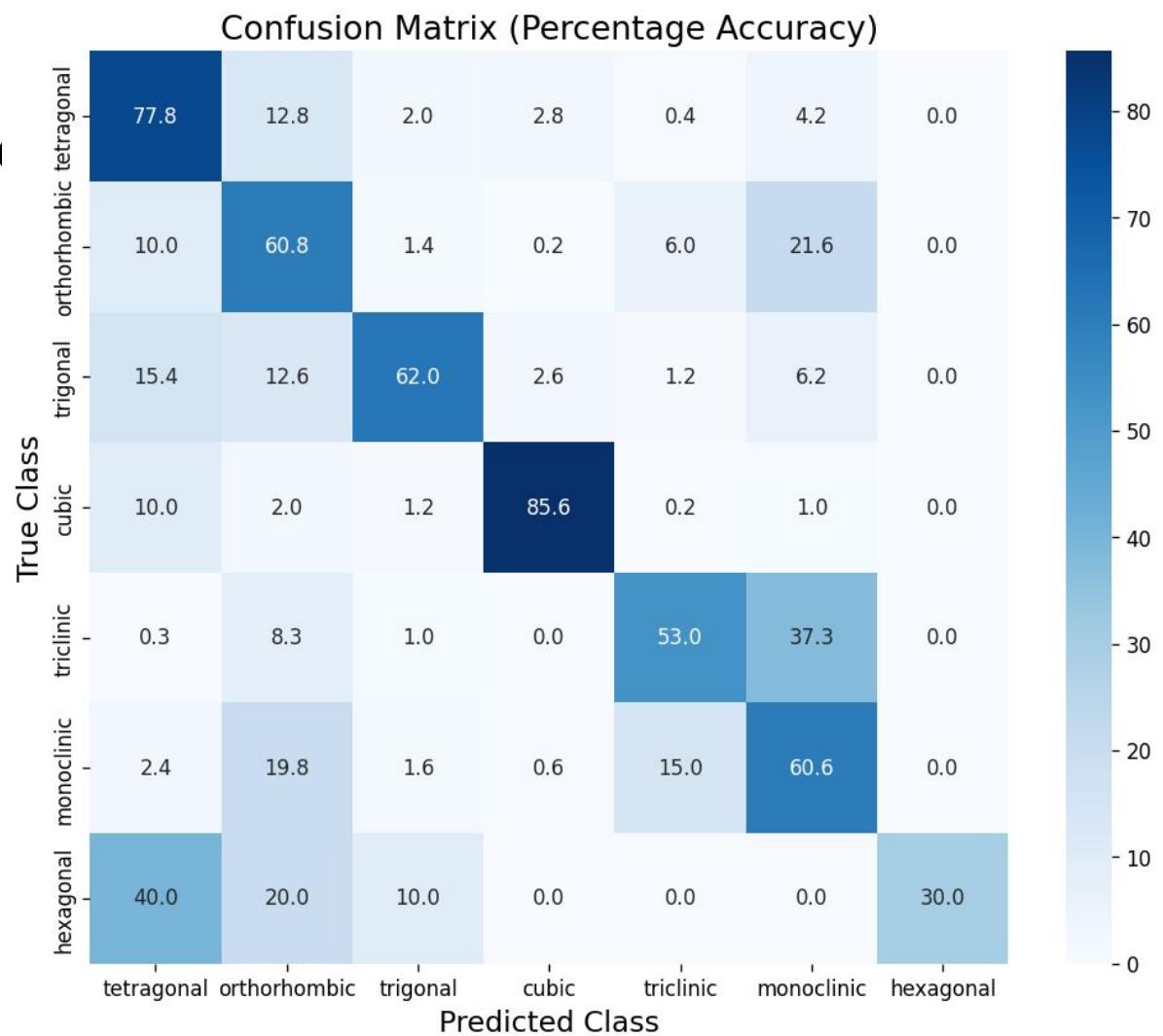
# Extra Trees Classifier - Peak



# Extra-Trees Accuracy

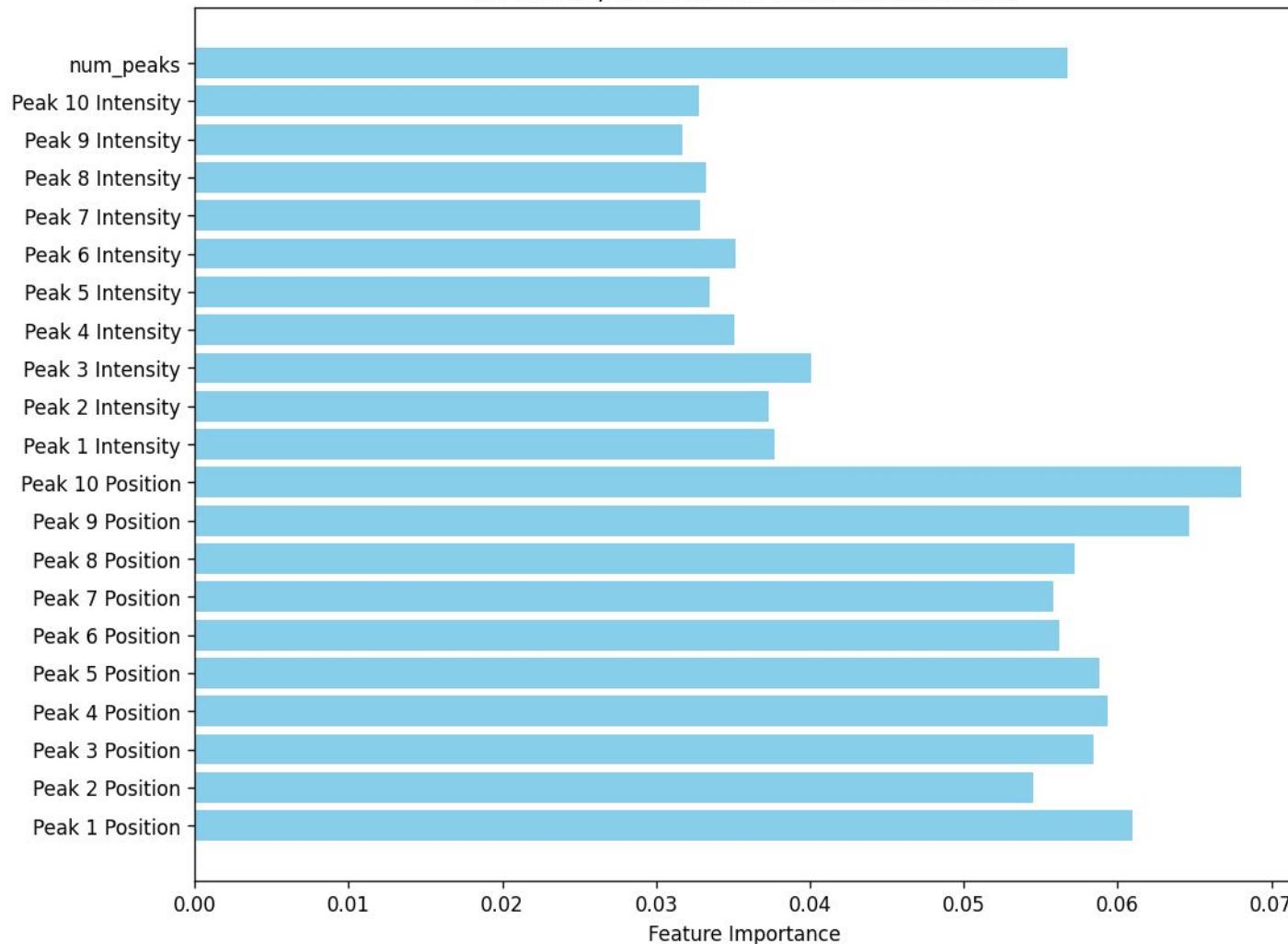
**Peak data 67.73 %**

# Extra Tree



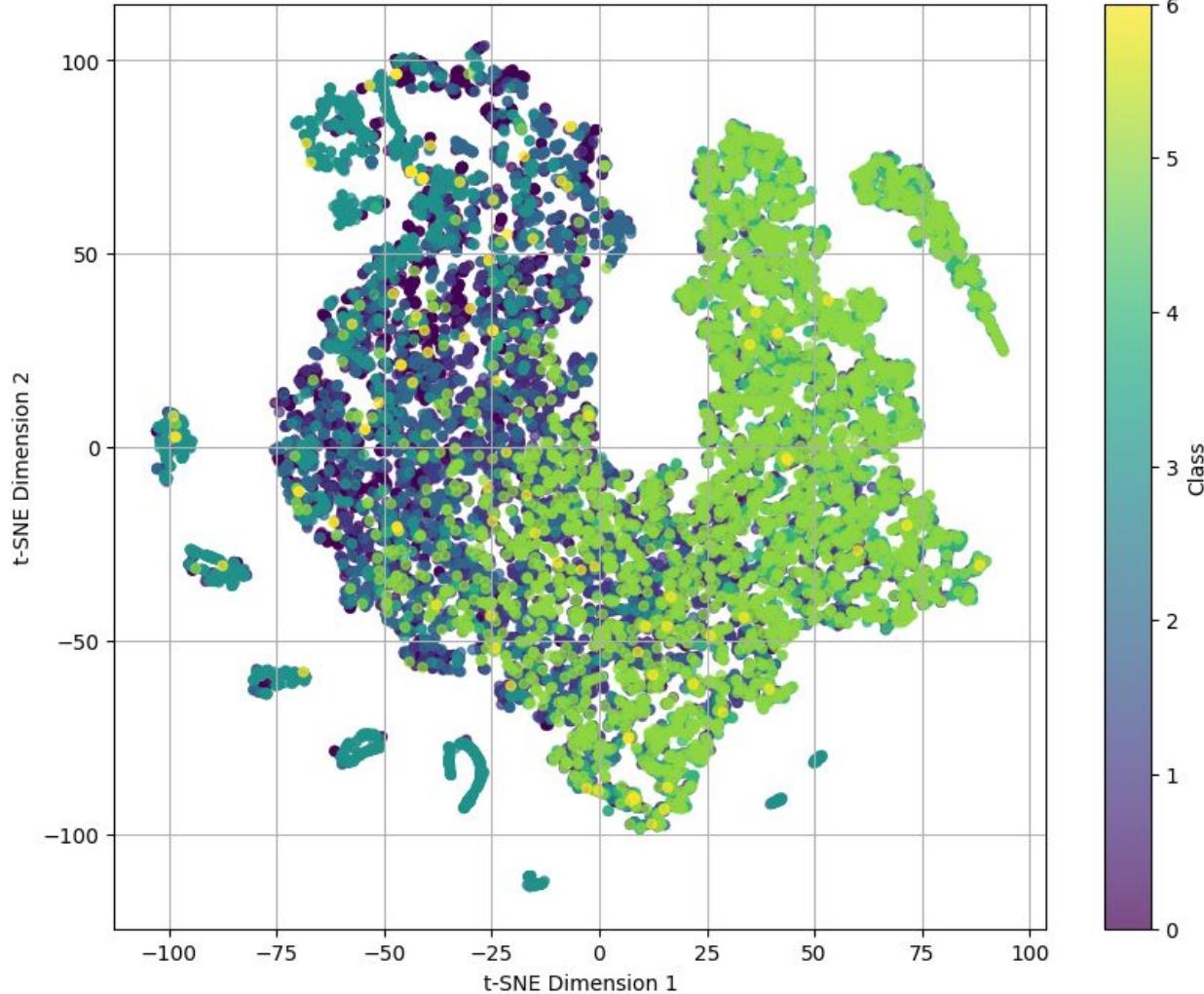
### Feature Importances from ExtraTreesClassifier

Extra

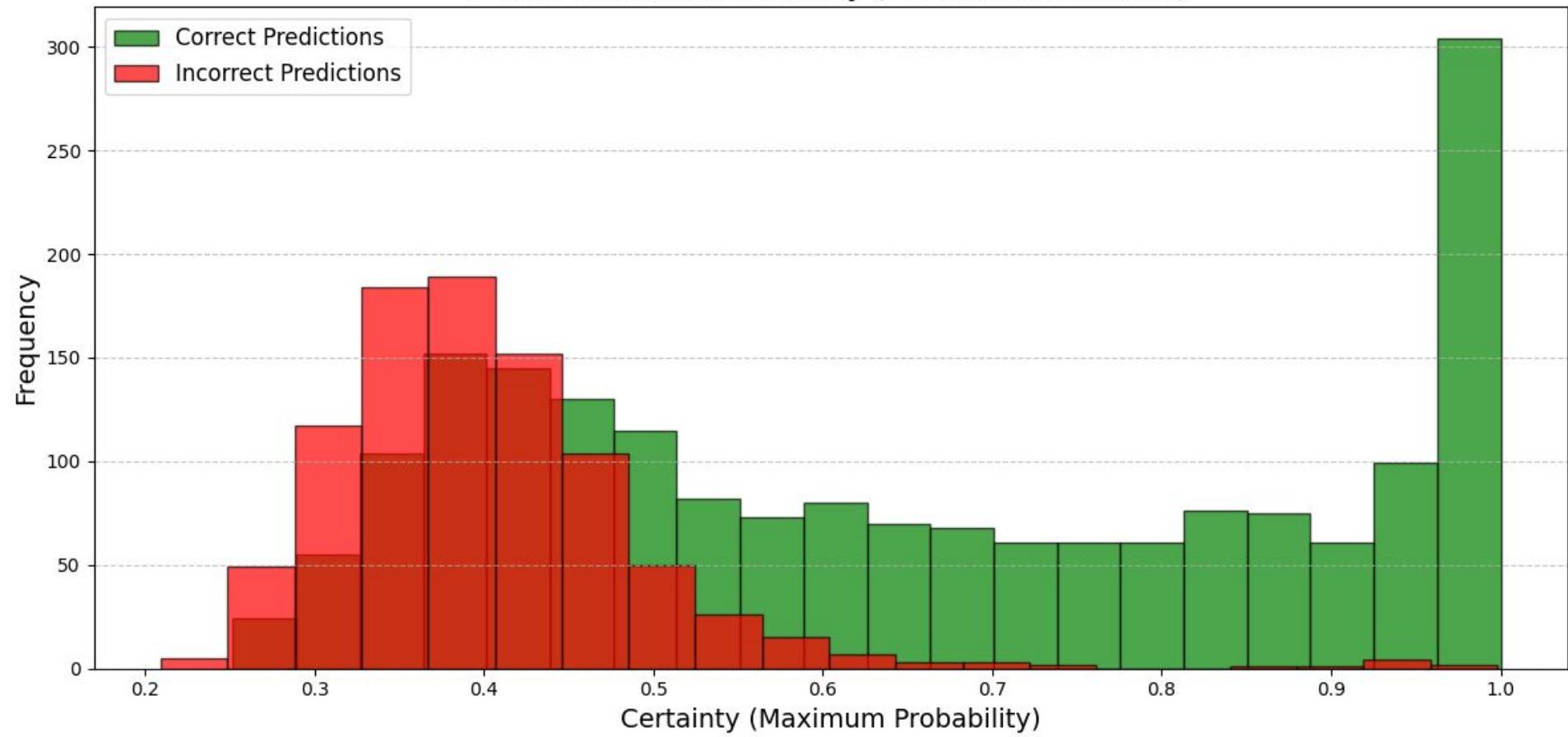


t-SNE Projection of Training Data

Extra Tree

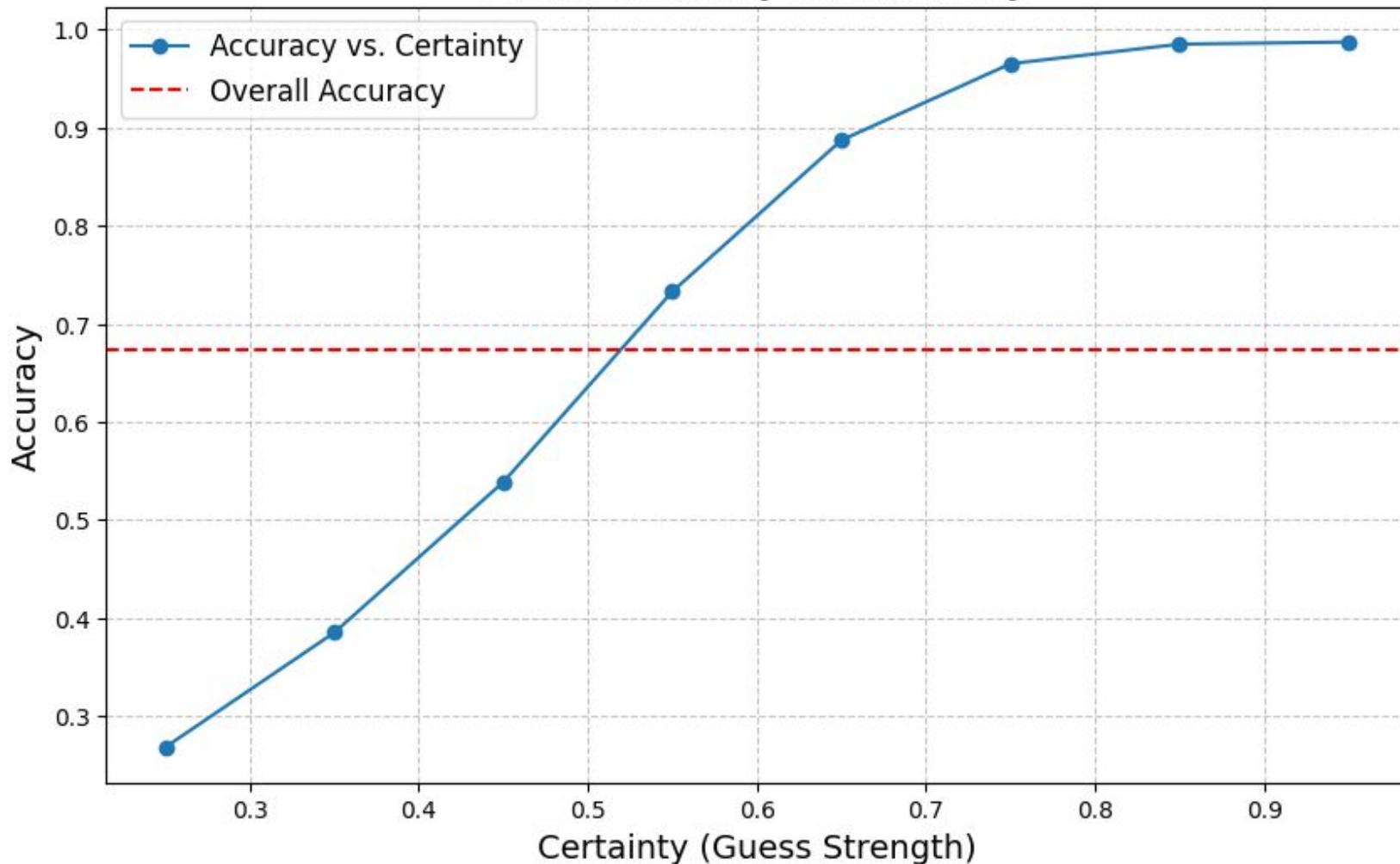


## Model Prediction Certainty (Correct vs Incorrect)



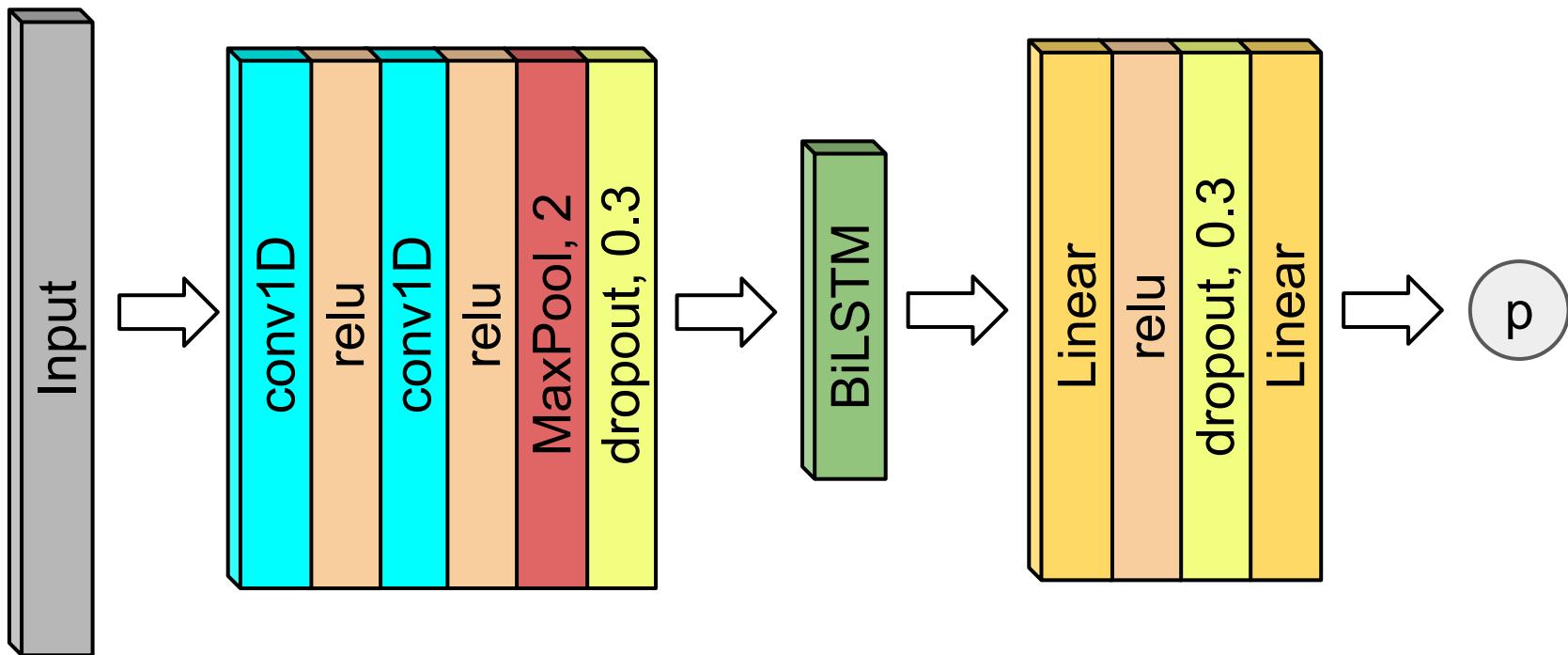
## Model Certainty vs. Accuracy

E

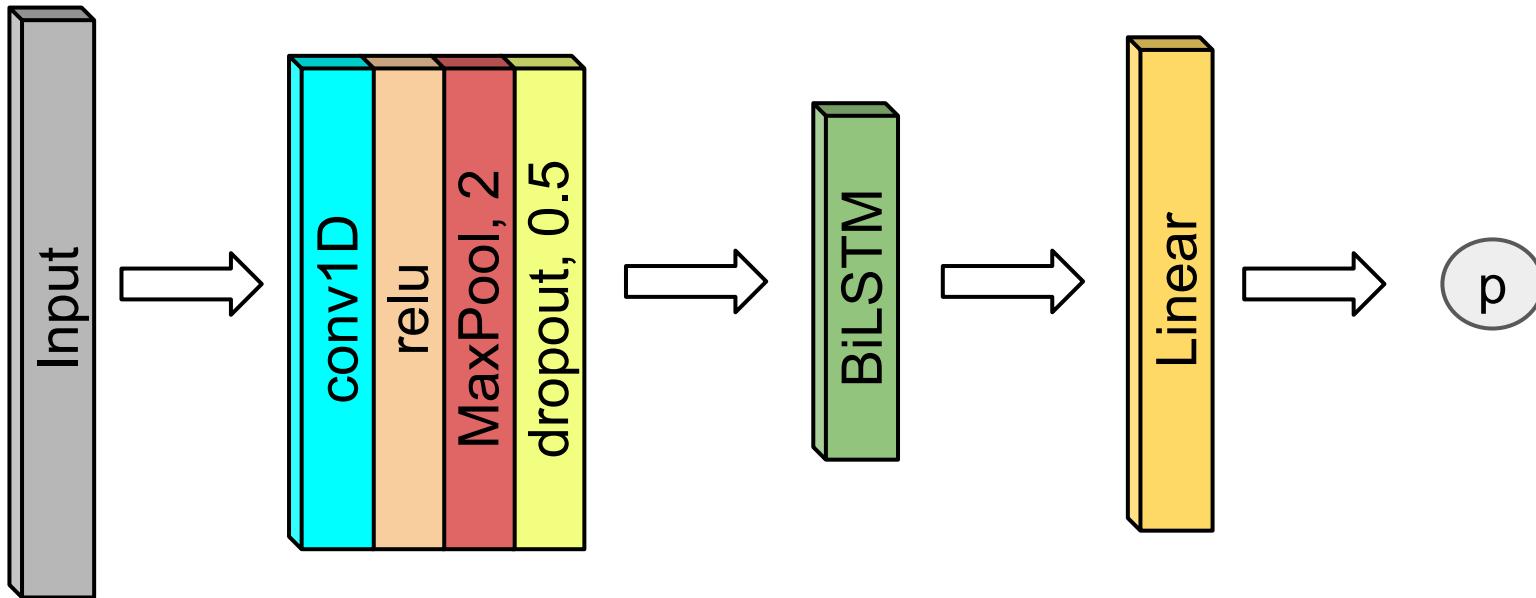


## 4. CNN + LSTM

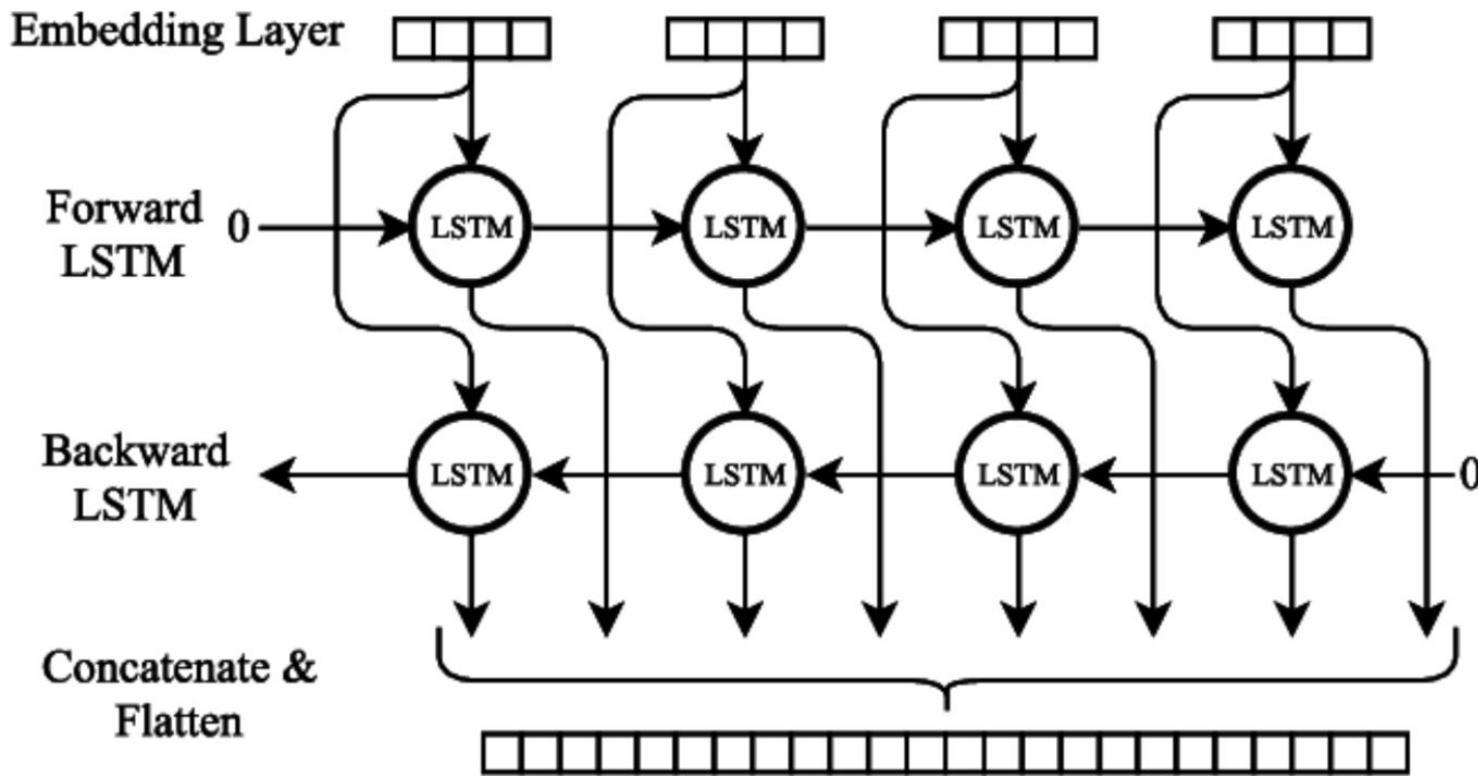
# CNN-LSTM: Architecture



# sCNN-LSTM: Architecture



# BiLSTM: Architecture



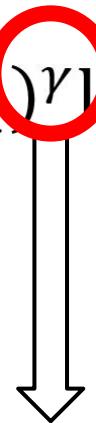
# Focal Loss

$$CE = -\log(p_t)$$

$$FL = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$

# Focal Loss

$$FL = -\alpha_t(1 - p_t)^\gamma \log(p_t)$$



weighting coefficient, set = 2

# Models' Accuracy

**All data:**

sCNN-LSTM, CE/ FL: 52.66%, 50.47%

**Peak data:**

CNN-LSTM, CE/ FL: **64.65%**, 54.88%

sCNN-LSTM, CE/ FL: 43.19%, 42.96%



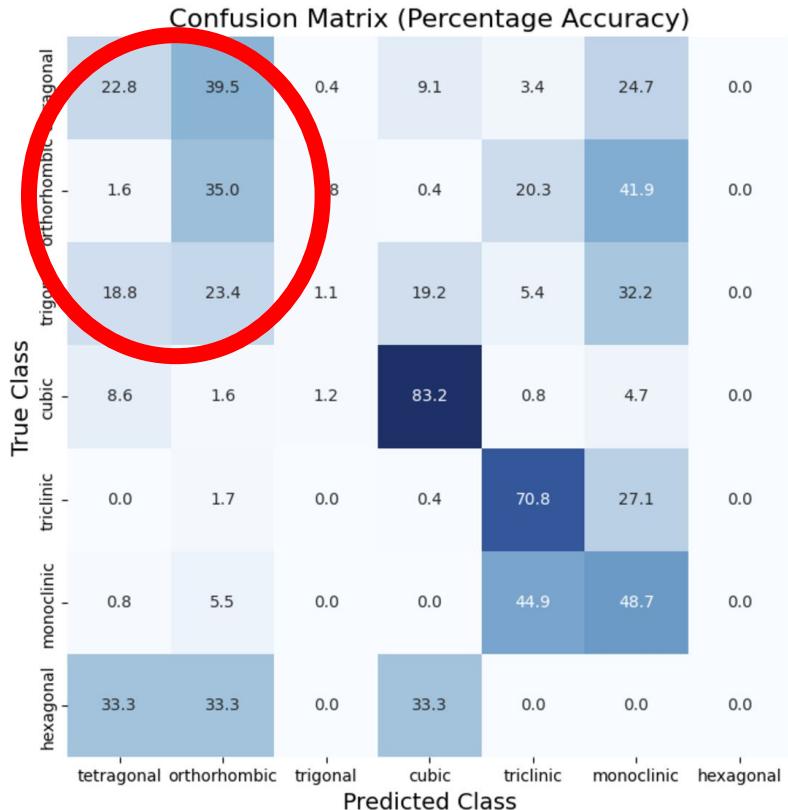
6 models

Comparison:

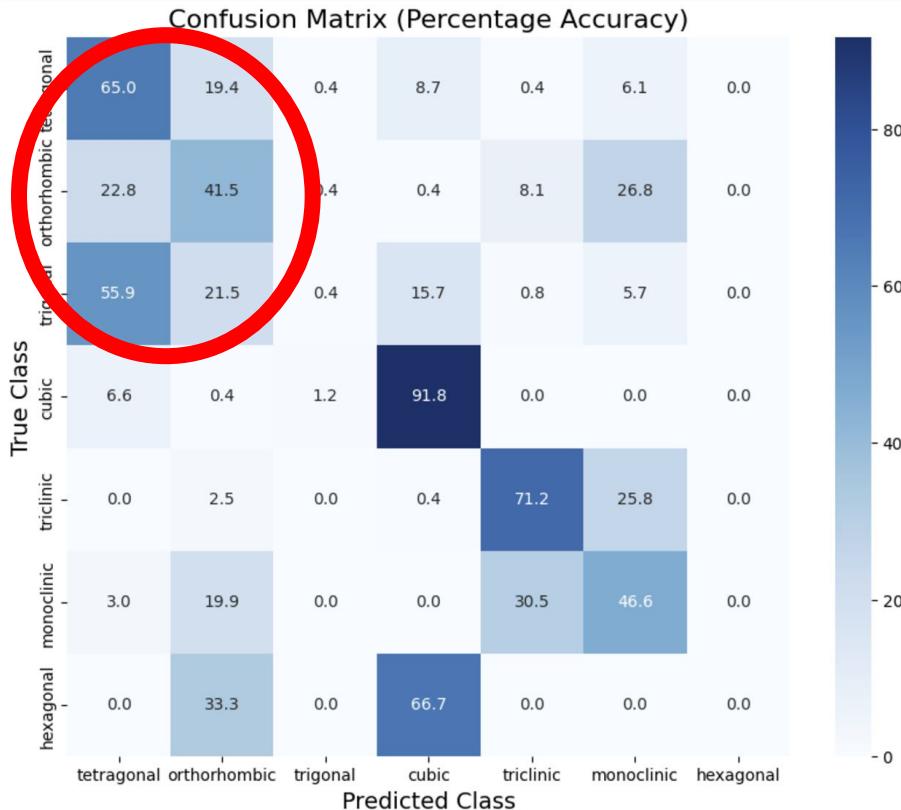
Peak Data 43.19% vs. All Data 52.66%

sCNN-LSTM, Cross Entropy

# Peak Data/ All Data - Confusion Matrix



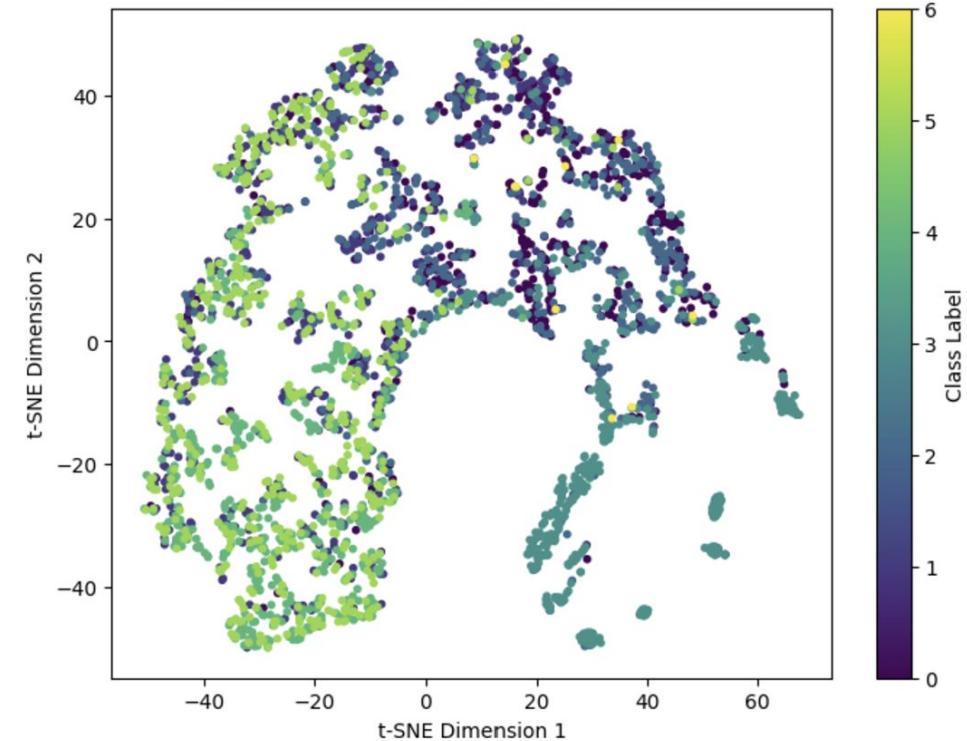
Peak Data



All Data

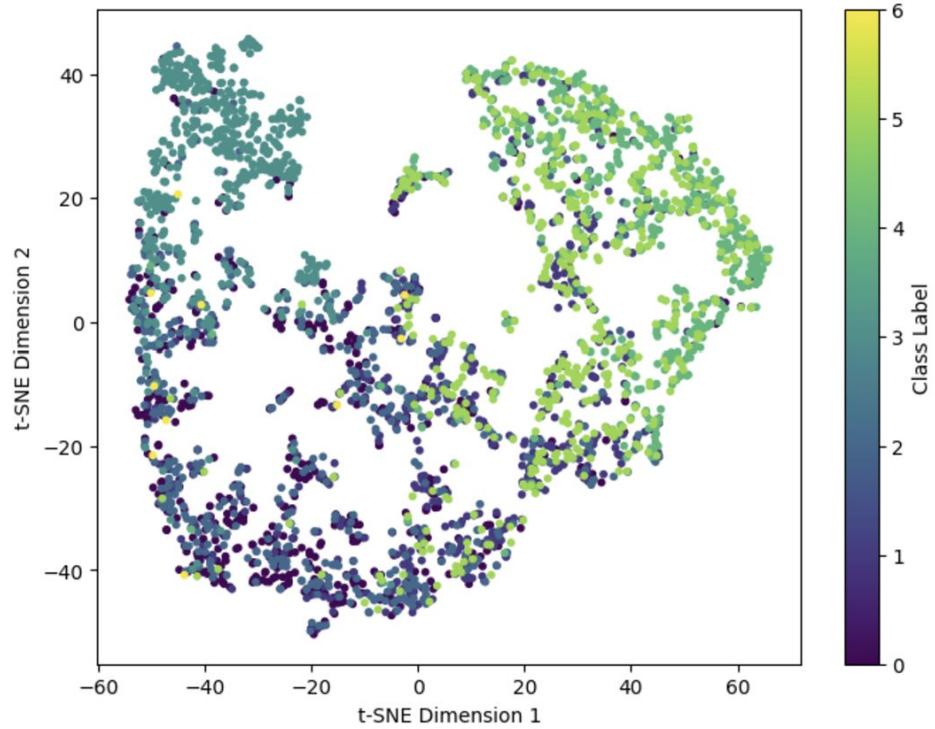
# Peak Data/ All Data - TSNE

t-SNE Visualization of XRD CNN Features



Peak Data

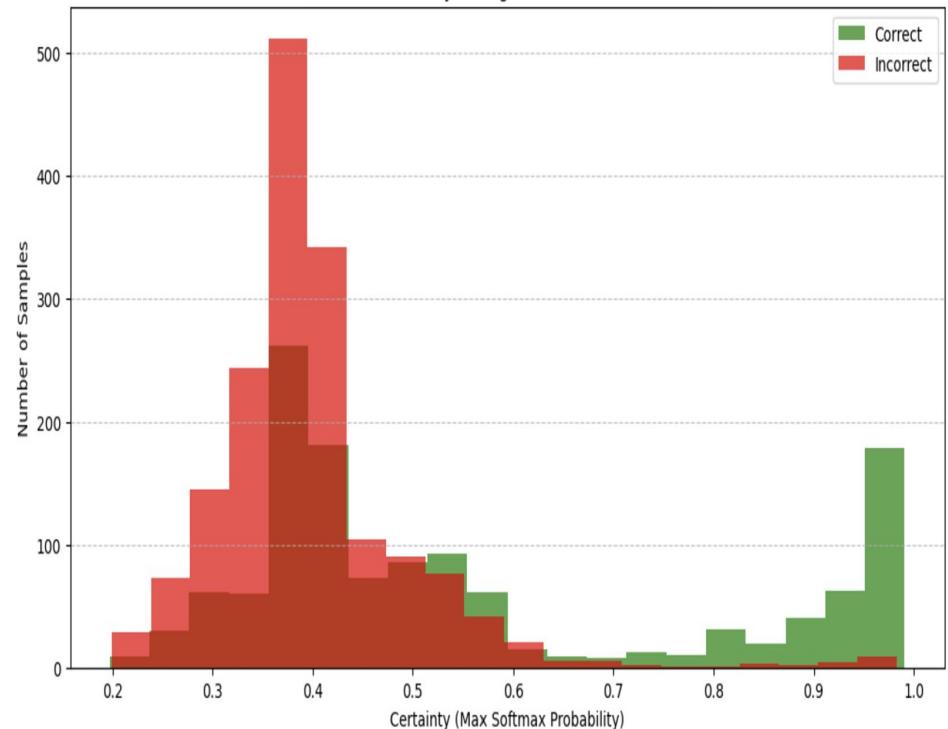
t-SNE Visualization of XRD CNN Features



All Data

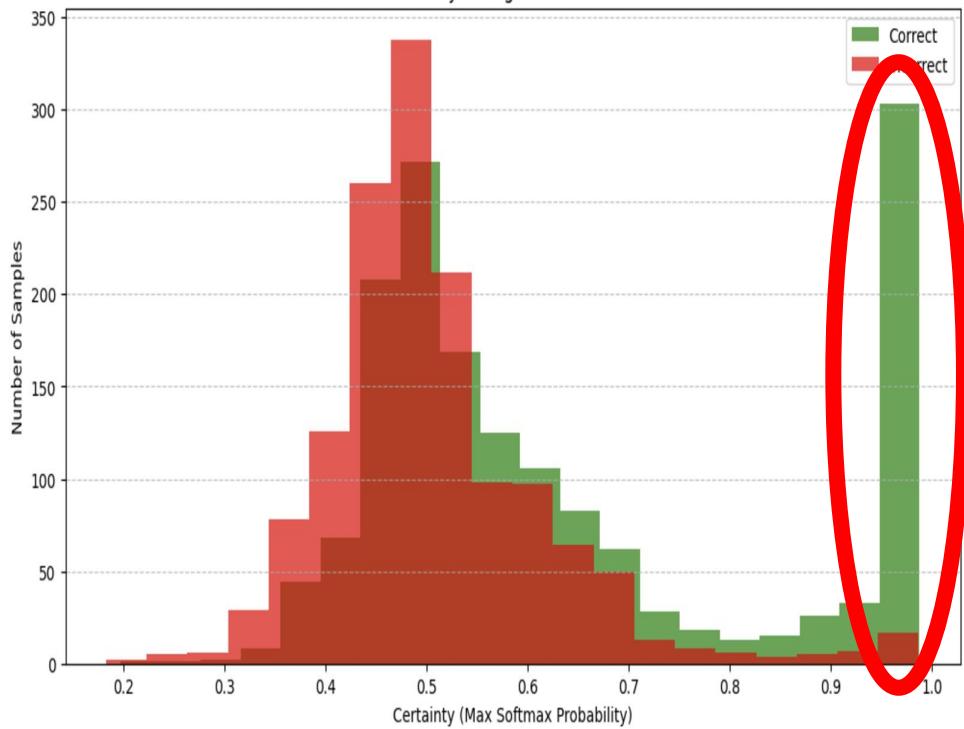
# Peak Data/ All Data - Network Certainty

Network Certainty Histogram: Correct vs Incorrect



Peak Data

Network Certainty Histogram: Correct vs Incorrect



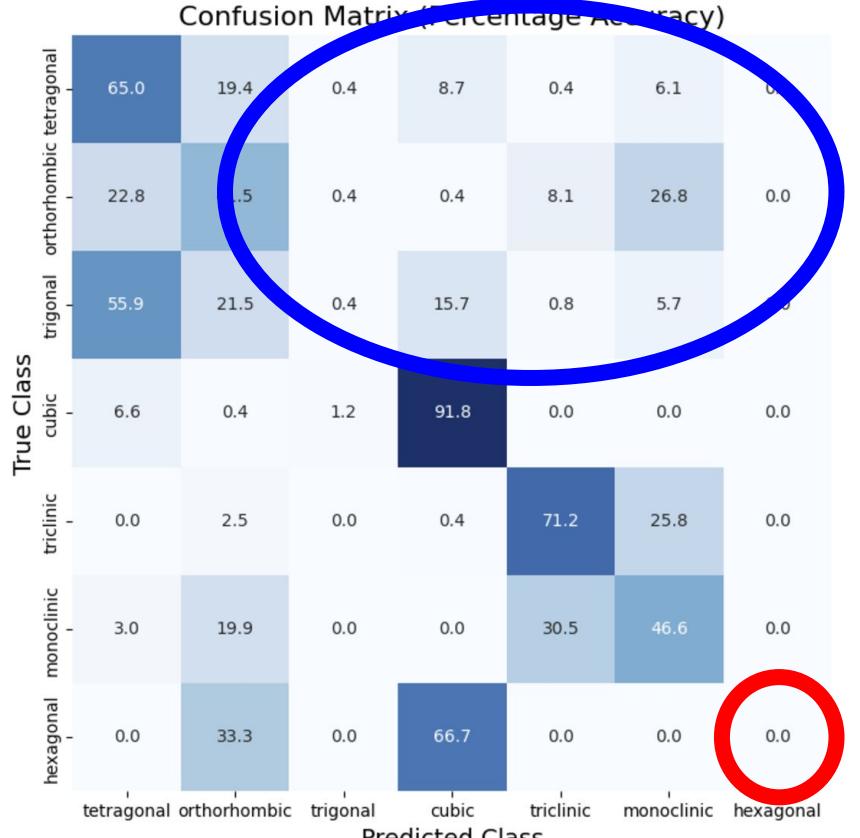
All Data

Comparison:

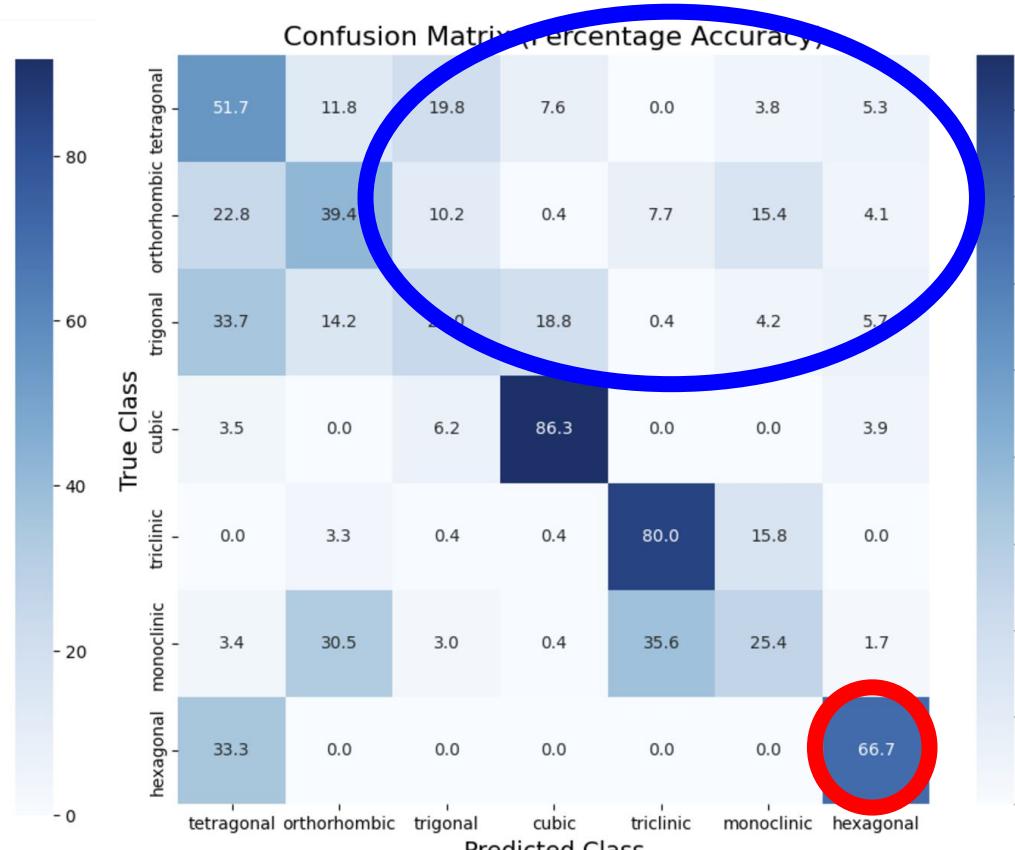
Cross Entropy 52.66% vs. Focal Loss 50.47%

sCNN-LSTM, All Data

# CE/ FL - Confusion Matrix

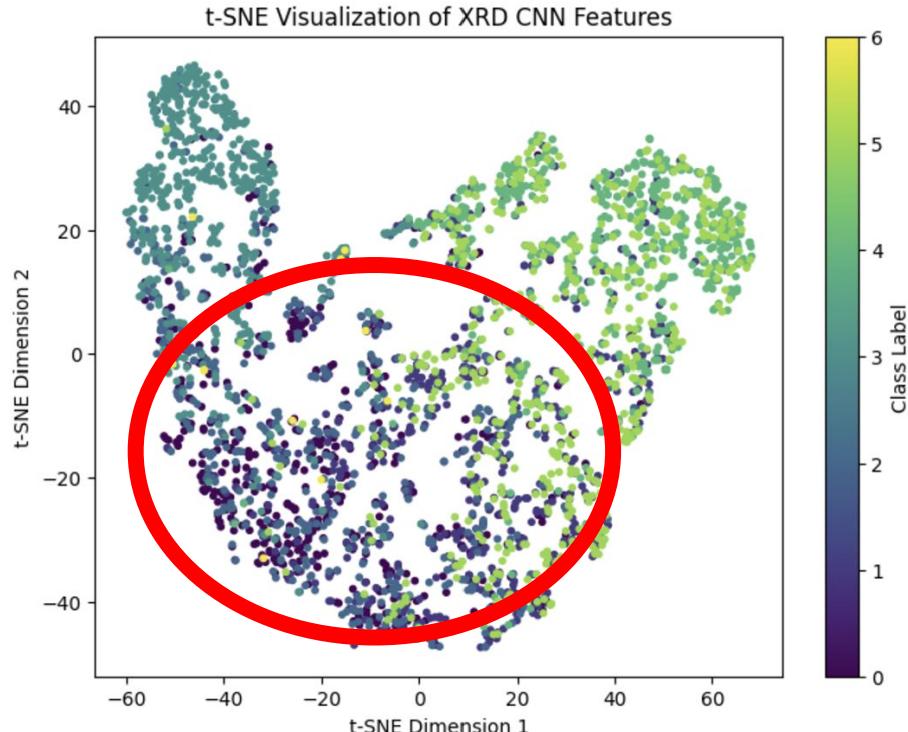
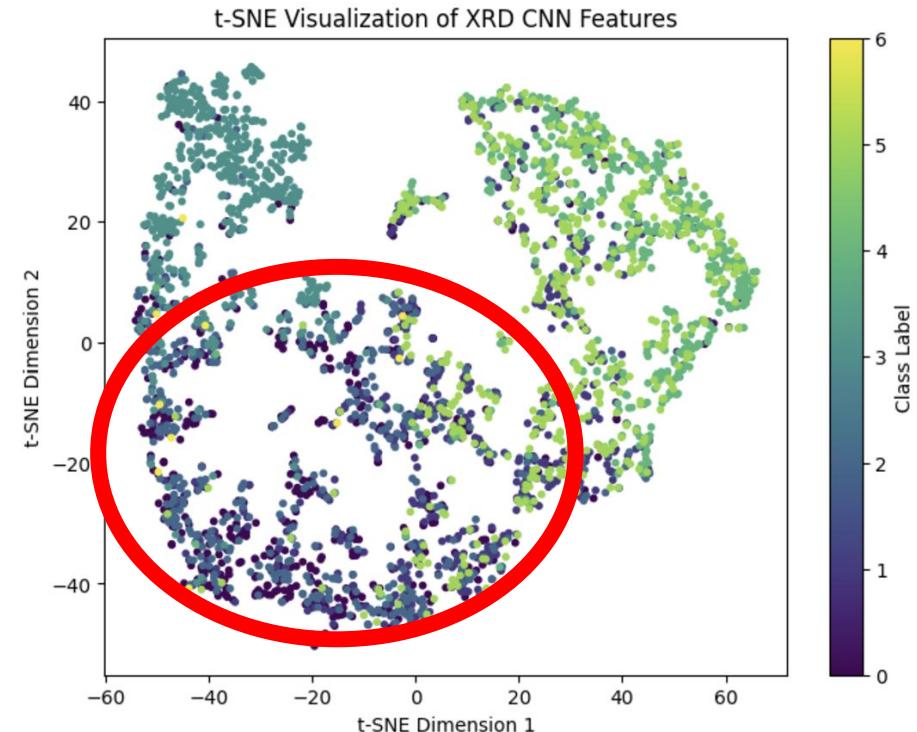


Cross Entropy

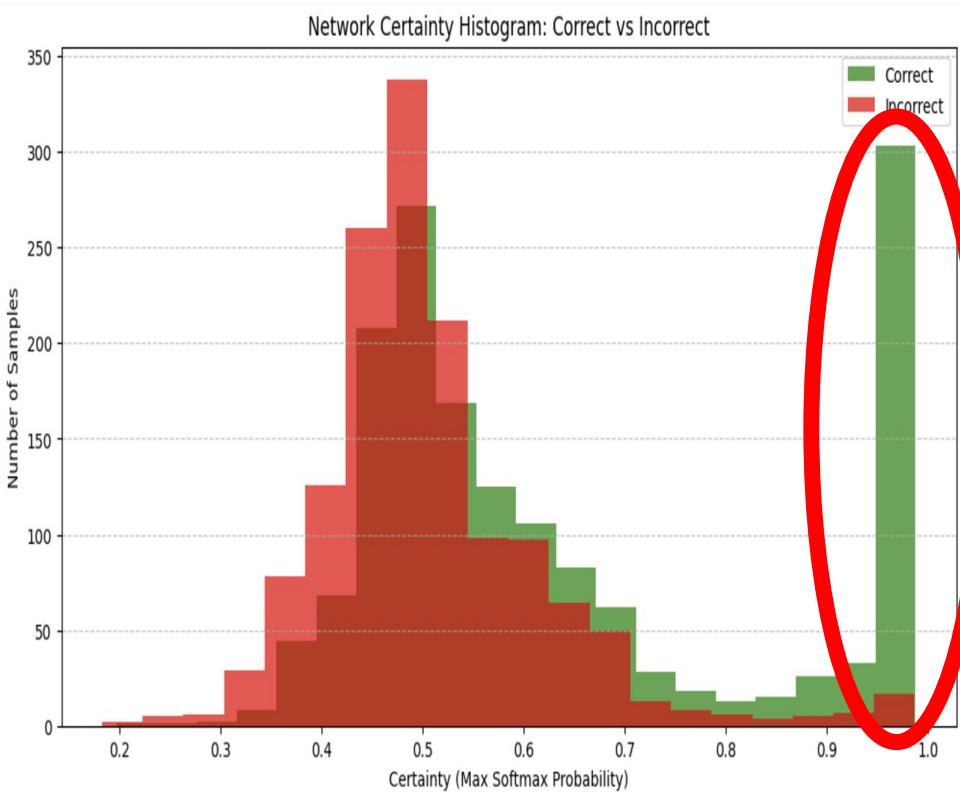


Focal Loss

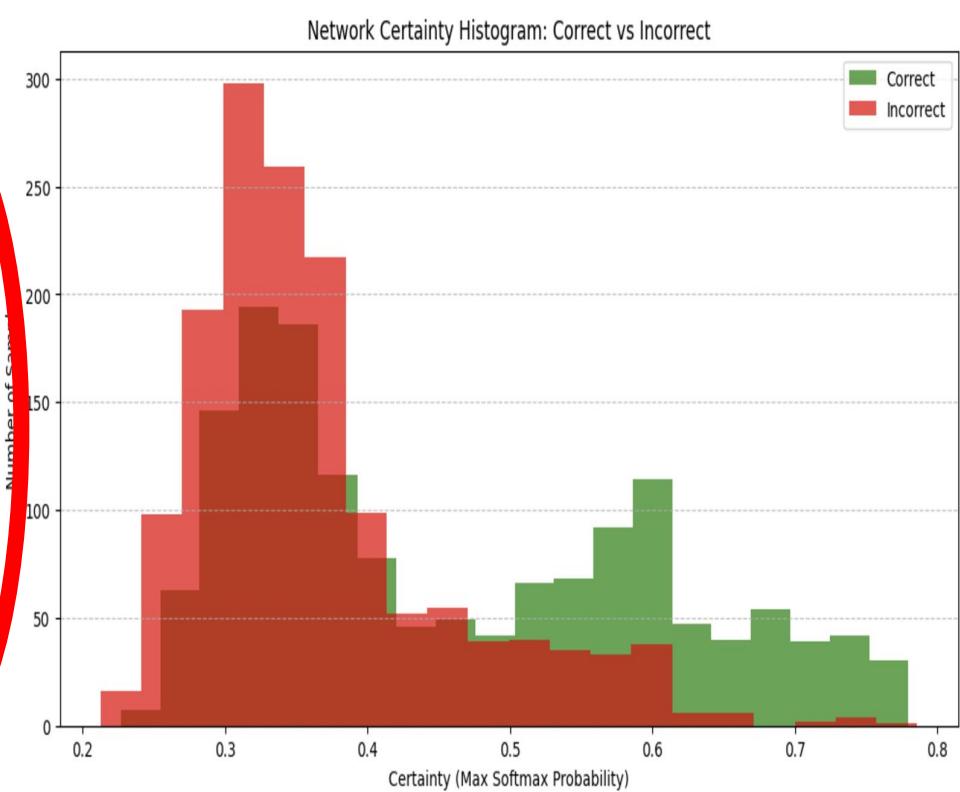
# CE/ FL - TSNE



# CE/ FL - Network Certainty



Cross Entropy



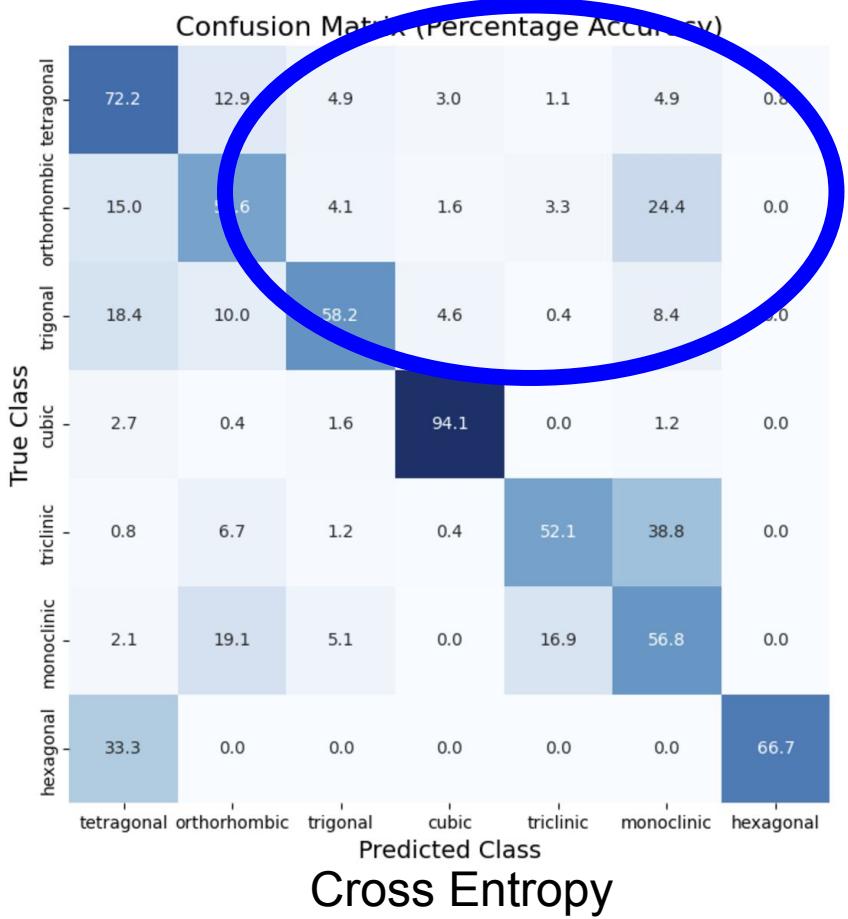
Focal Loss

Comparison:

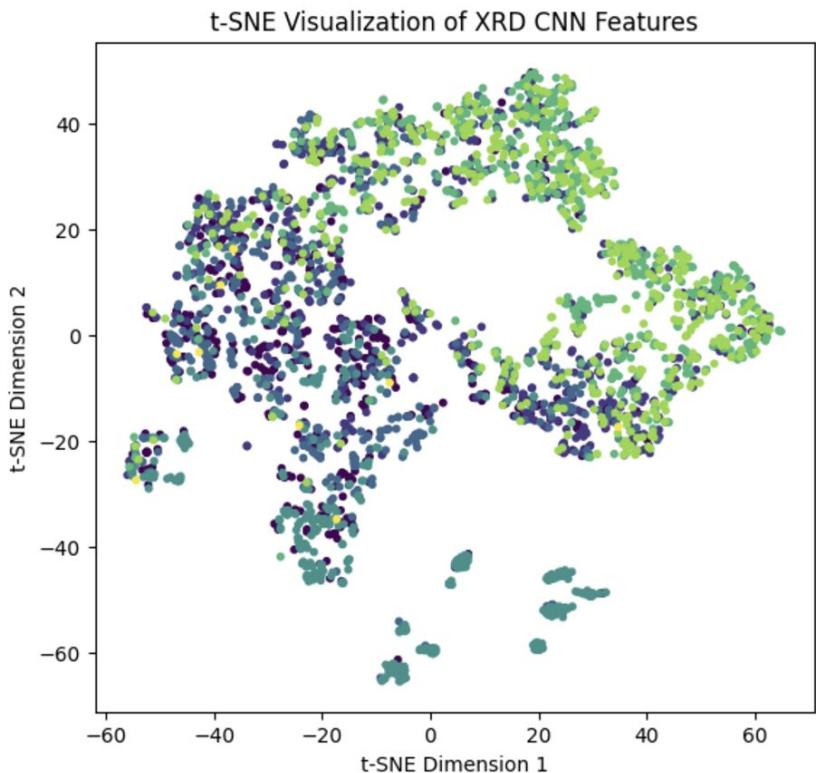
Cross Entropy 64.65% vs. Focal Loss 54.88%

CNN-LSTM, Peak Data

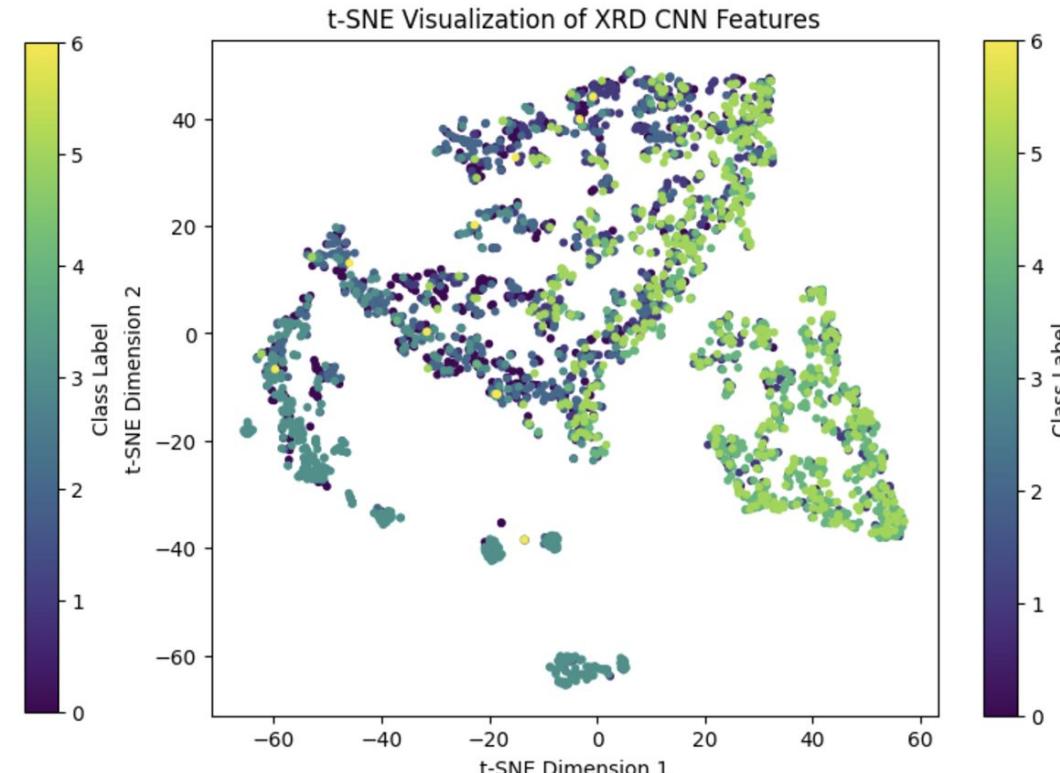
# CE/ FL - Confusion Matrix



# CE/ FL - TSNE

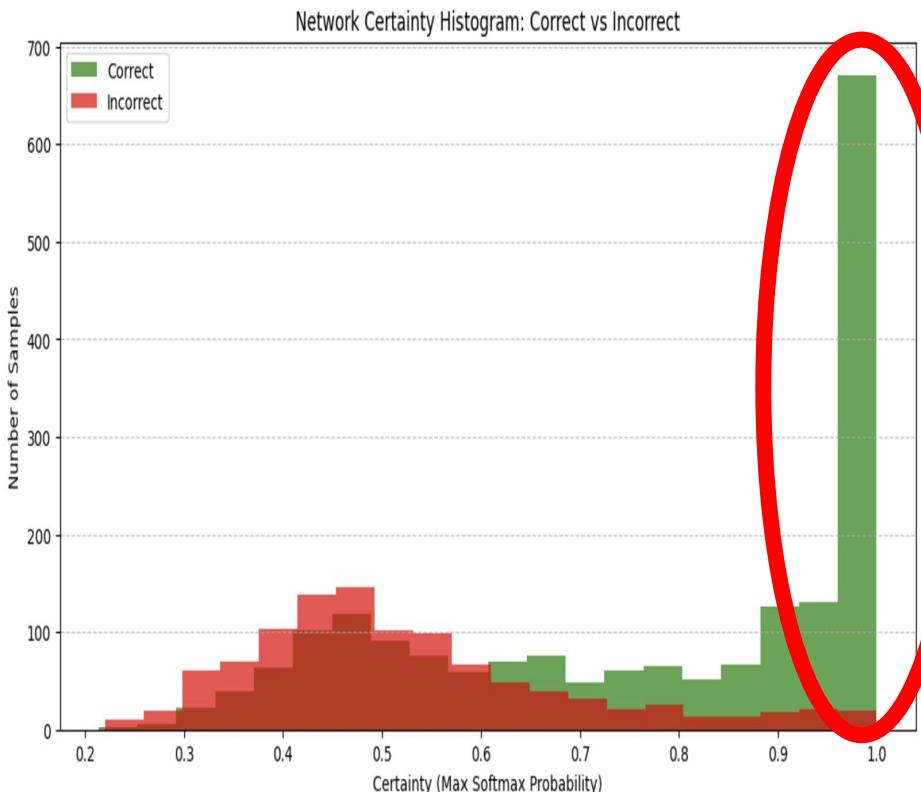


Cross Entropy

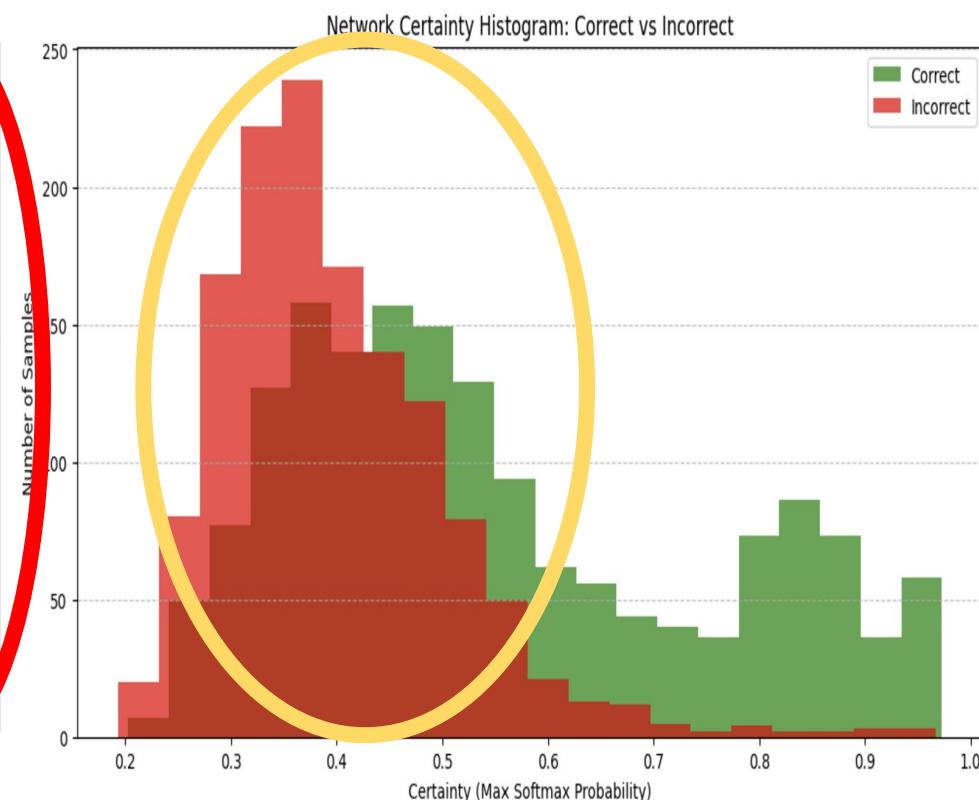


Focal Loss

# CE/ FL - Network Certainty



Cross Entropy

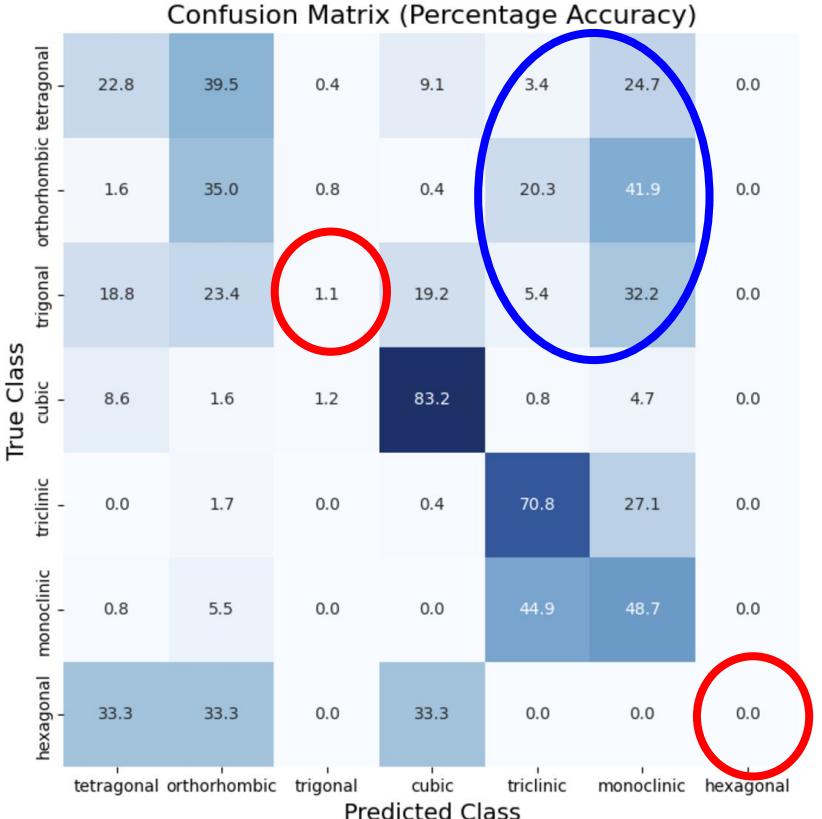


Focal Loss

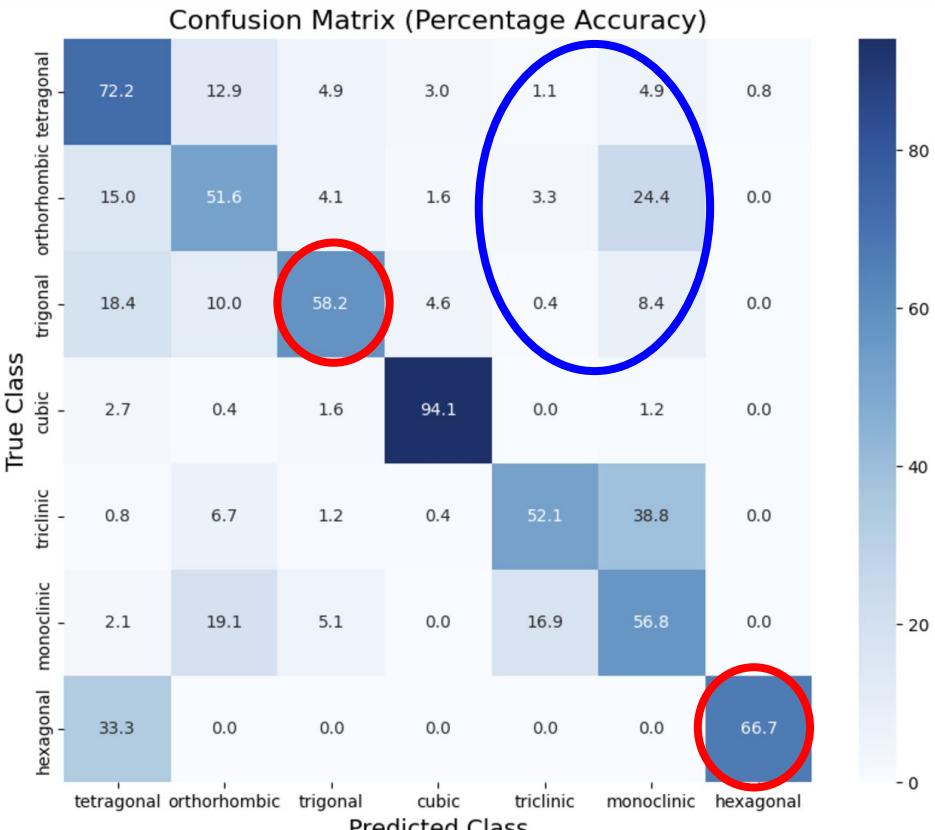
Comparison:  
Small Model 43.19% vs. Big Model 64.65%

Cross Entropy, Peak Data

# Peak Data - Confusion Matrix

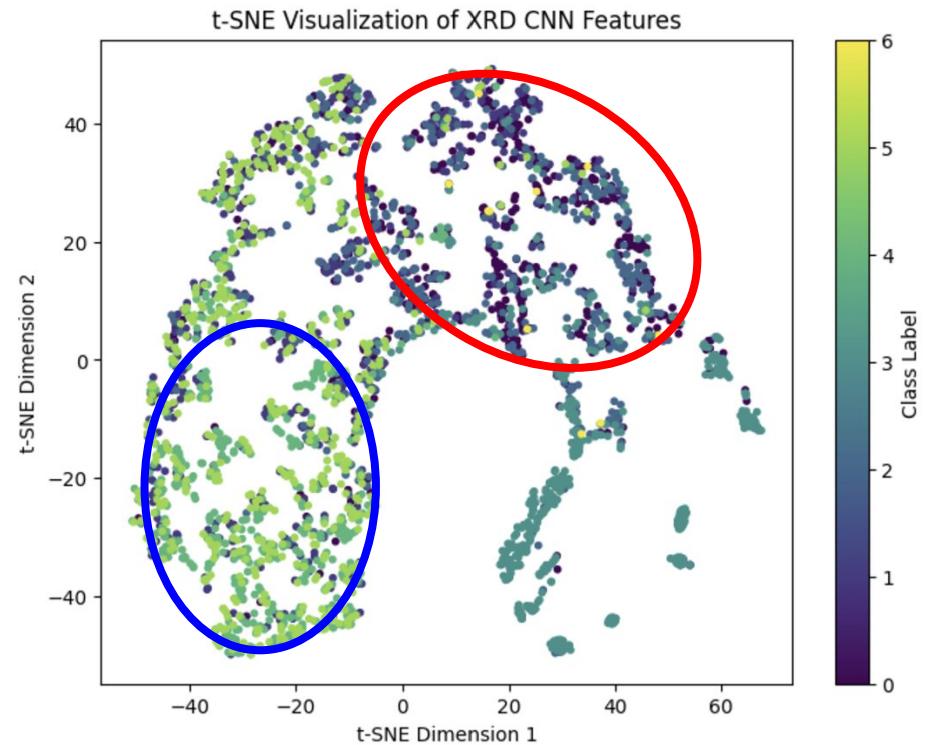


SCNN-LSTM

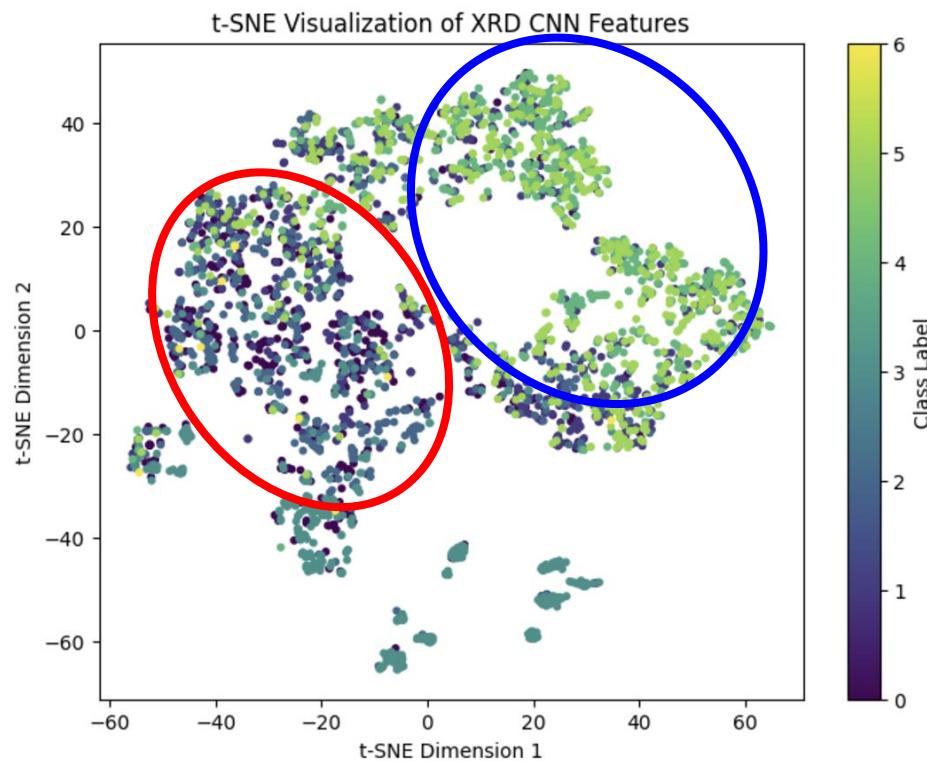


CNN-LSTM

# Peak Data - TSNE



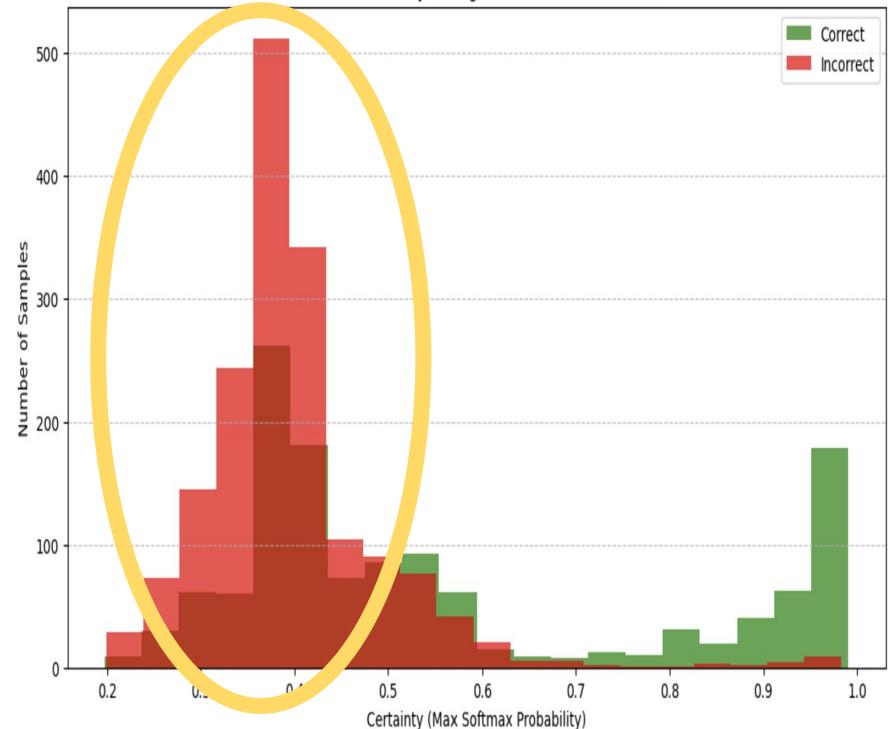
SCNN-LSTM



CNN-LSTM

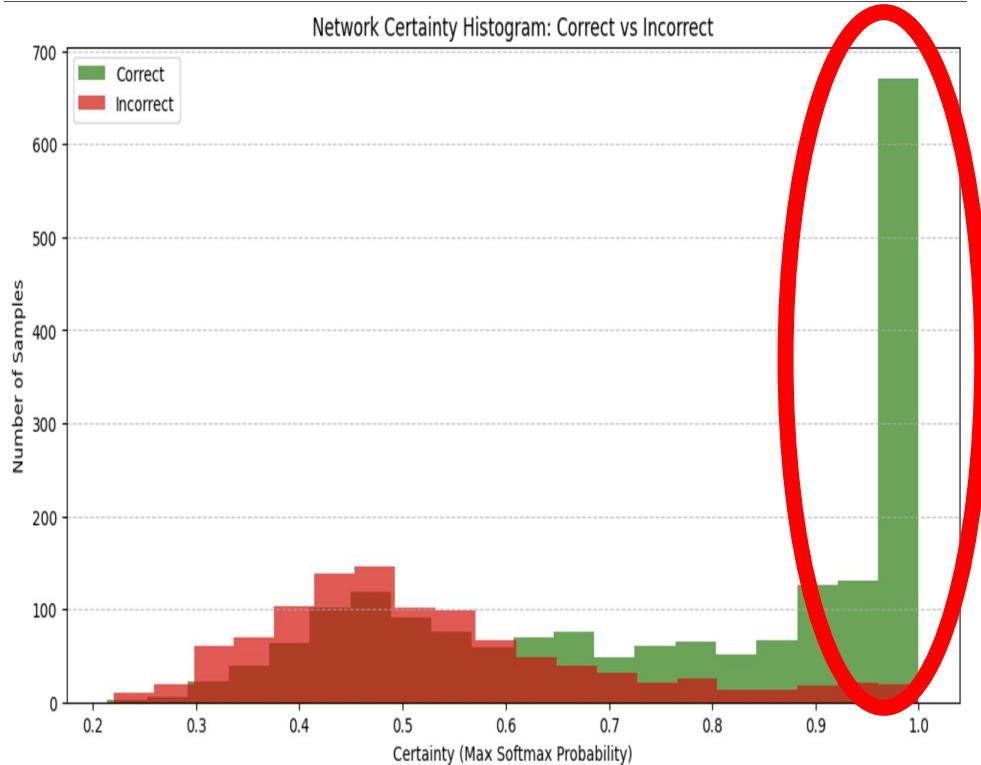
# Peak Data - Network Certainty

Network Certainty Histogram: Correct vs Incorrect



SCNN-LSTM

Network Certainty Histogram: Correct vs Incorrect



CNN-LSTM

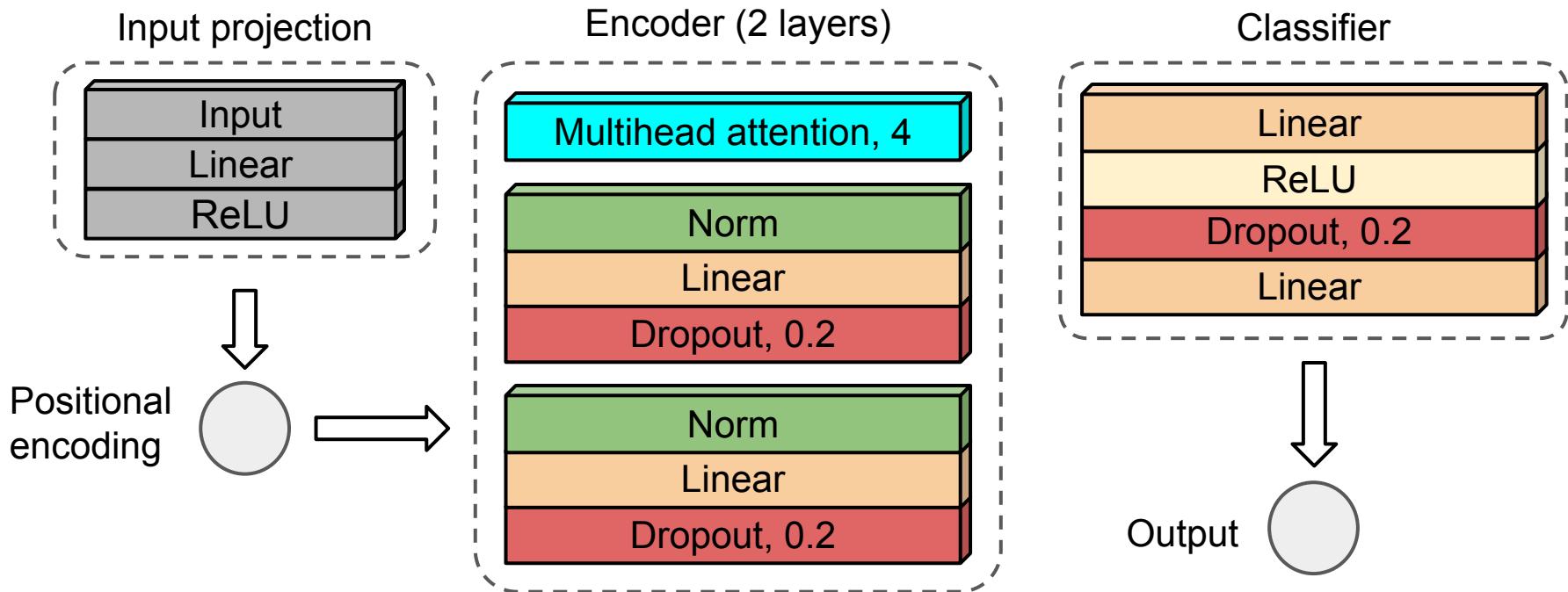
## Small Conclusion

1. sCNN-LSTM works better on big dataset.
2. Focal loss does well on small samples but sacrifices overall performance.
3. CNN-LSTM (big model) works better on peak dataset than sCNN-LSTM.

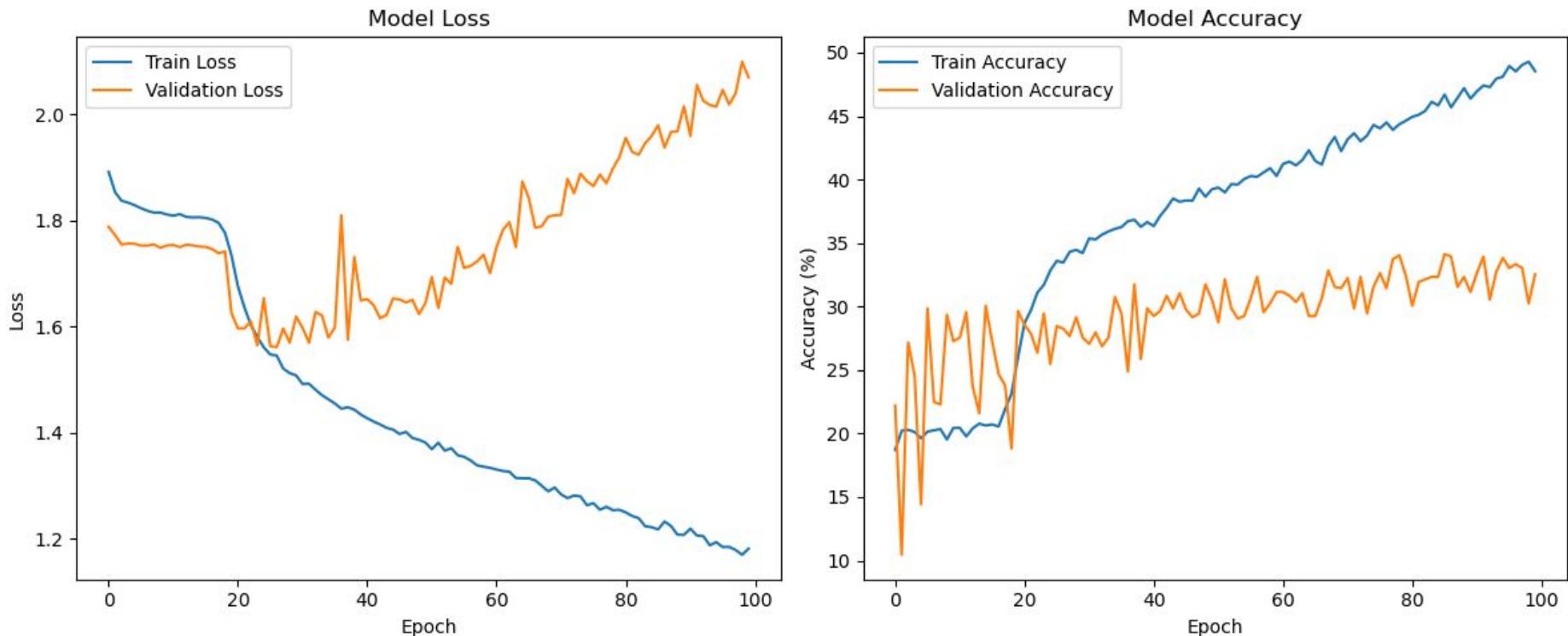
# 5. Transformer



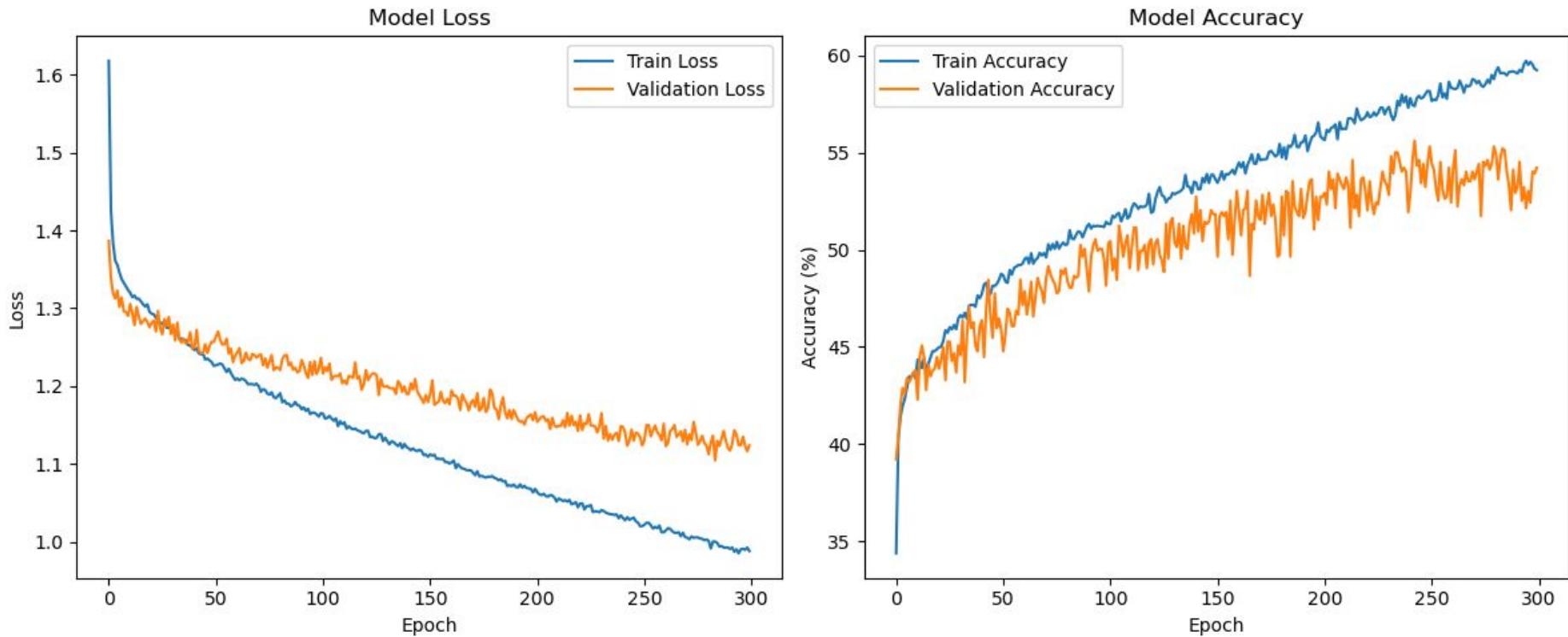
# Architecture



# Full spectrum

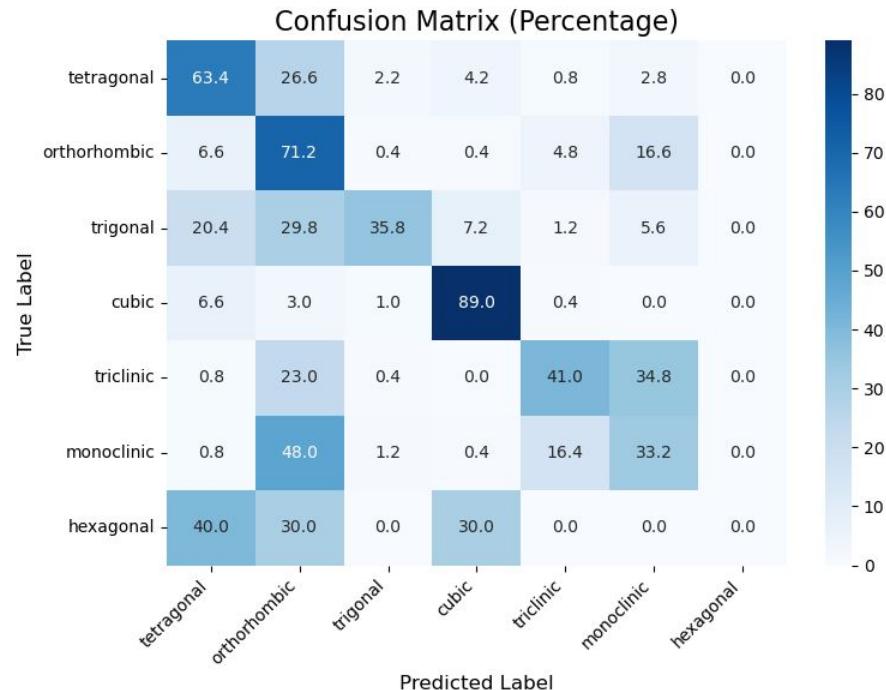
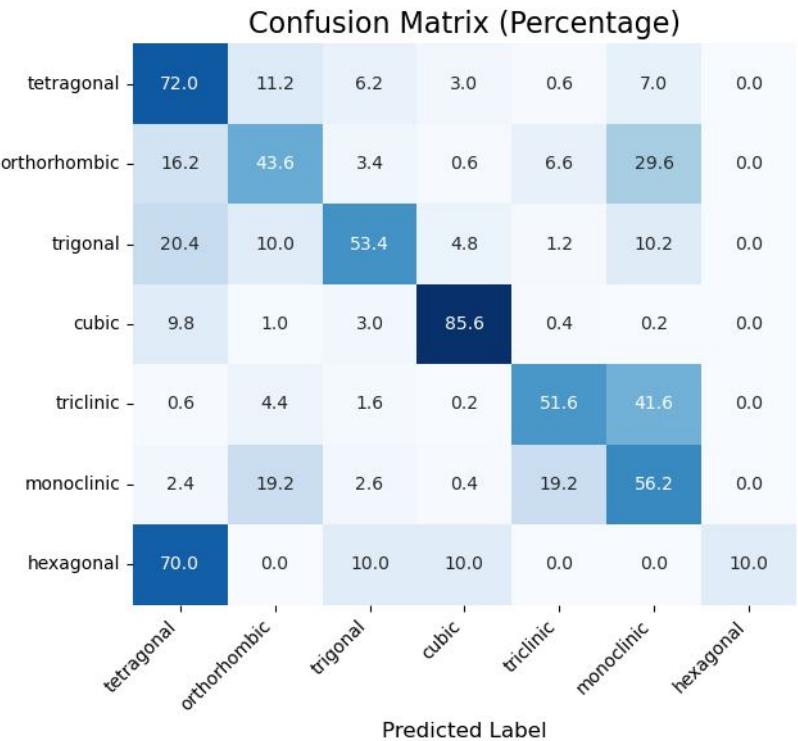


# Peaks only

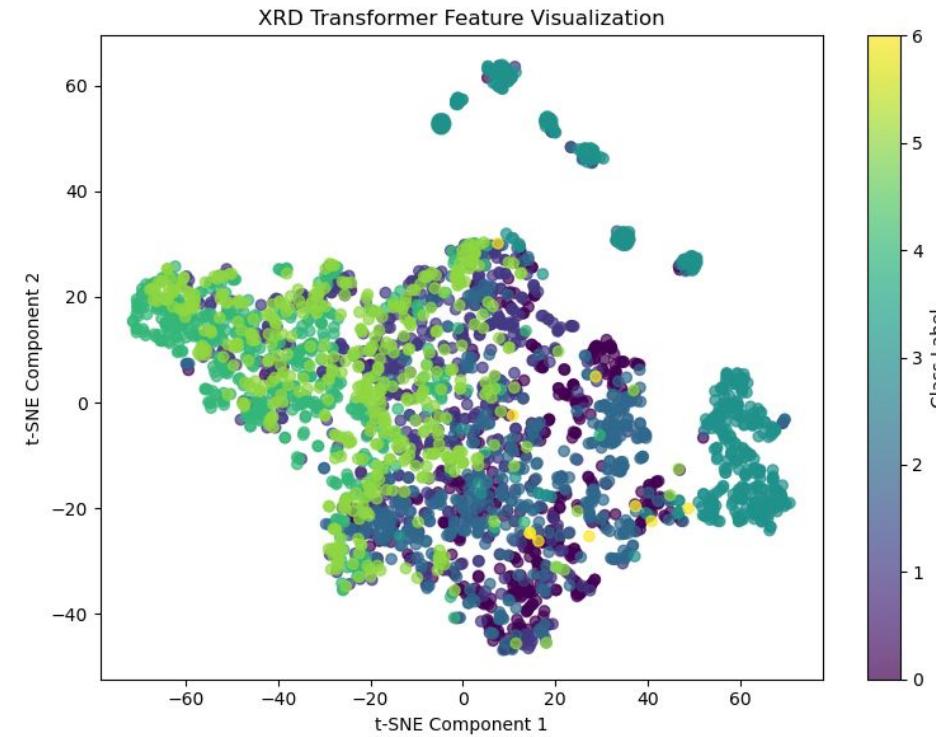
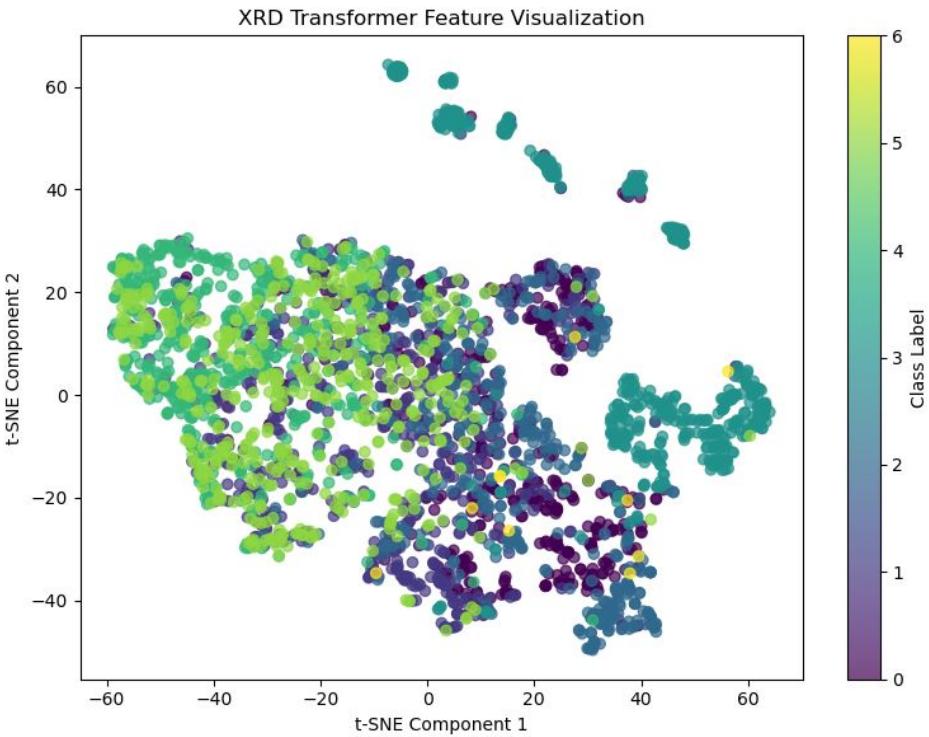


# Confusion matrices (best acc: 60%)

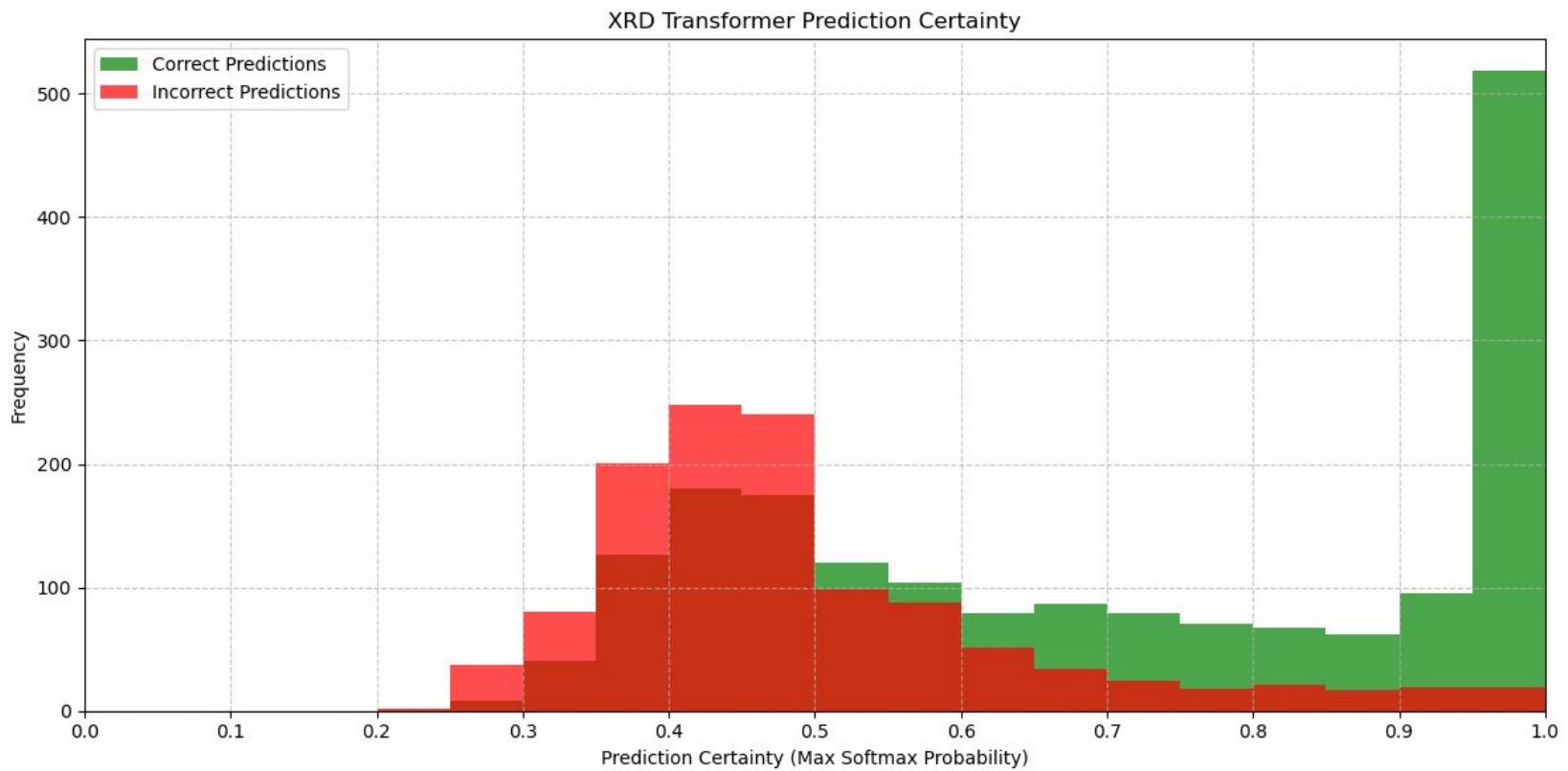
True Label



# Transformer t-SNE



# Certainty



# Summary / Final remarks

# Our models

Model	Full spectrum (%)	Peaks only (%)
Dense (150 MB)	35	56
CNN+Dense (160 MB)		60
CNN-LSTM with CE/FL (1 MB)		64/54
sCNN-LSTM with CE/FL (75 KB)	52/50	43/42
Extreme Randomized Tree (90 MB)		67
Transformer (460 KB)	33	60

# Future improvements

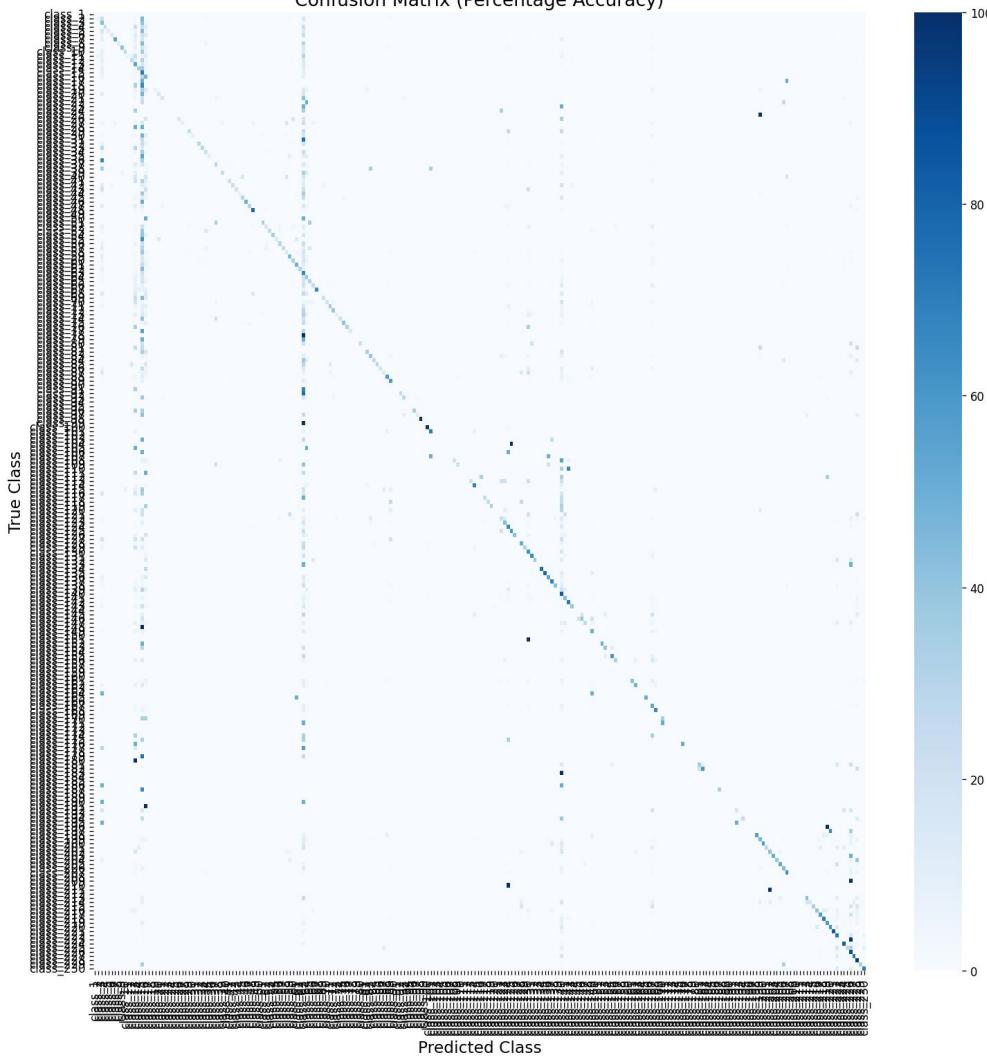
## Training-wise

- Extract other physical features
- Try new models (e.g., RNN)

## Data-wise

- Larger datasets
- Filter noise
- Normalized spectrum

Confusion Matrix (Percentage Accuracy)



# Supplements

We have trained the  
Classifier.

Random Extra Trees

# References

[1] PhysRevB.99.245120

[2] Nature Scientific Reports “Symmetry prediction and knowledge discovery from X-ray diffraction patterns using an interpretable machine learning approach”