

### MOTIVATION

Our motivation is to classify facial expressions using different machine learning models. We are using the FER (Facial Expression Recognition) - 2013 dataset by Manas Sambare.

### AIM

Classifying facial expressions using Convolutional Neural Networks and Vision Transformers

### METHODS

**FERNet:** A CNN model with architecture inspired by AlexNet. Initial hyperparameters were chosen based on previous projects on Kaggle. A random hyperparameter search was deployed to improve the model's test accuracy. The tuned model is trained until convergence.

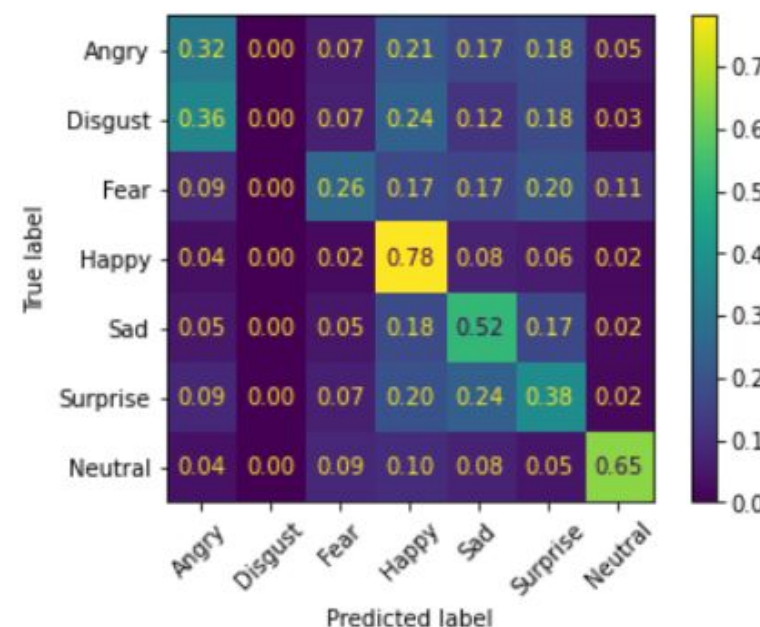
**ResNet:** A Residual Network model is a variant of CNN models. It incorporates the concept of residual learning which solves many performance problems seen in classic CNN models. The main difference in the implementation of a CNN model is the use of residual blocks. Instead of stacking many nonlinear layers on top of each other, we incorporate an underlying identity mapping in the ResNet model between our stacked layers. We essentially give the mapping a reference point which in theory is easier to train from than from an unreferenced mapping. We ran the training on the same training and testing dataset as our AlexNet model.

**ViT Image Classifier:** Our model was developed using the HuggingFace Transformers library. The FER dataset was prepared using the pretrained "vit-base-patch16-224-in21k" feature extractor which is pre-trained on ImageNet-21K at resolution 224x224. The training dataset was scaled down to one-sixth of the original dataset but was made sure to have the same distribution of labels. The same pretrained model was also utilized by our ViT image classification model. The model was trained until convergence.

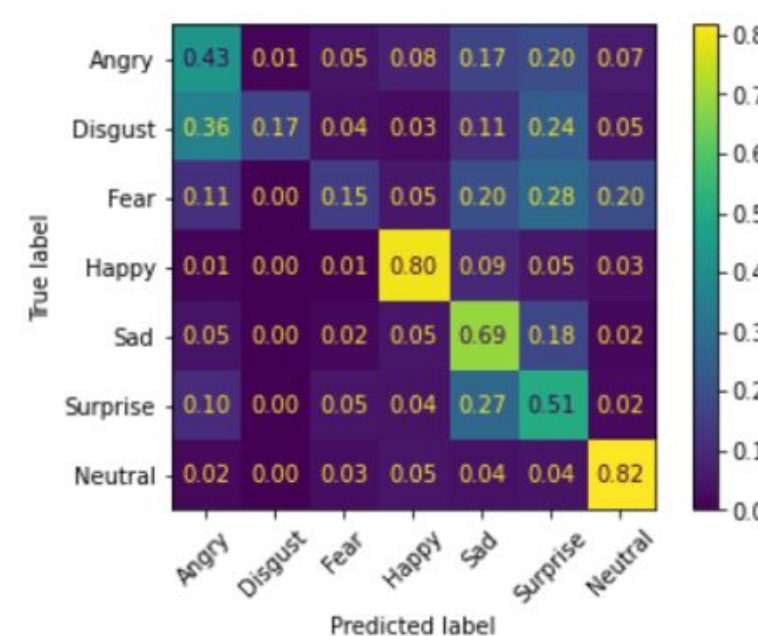
**Model Evaluation:** By keeping the dataset's characteristics consistent, we developed several tests to measure the performances of each model. These tests include the test accuracy, training-to-testing accuracy ratio, confusion matrix, and convergence rate. These scores were assessed to quantify the model's performance as well as identify any overfit.

**Literature Evaluation:** The results of the experiment are compared with results from published literature as well as from previous submissions on Kaggle. The process is carried out to verify the validity of our findings.

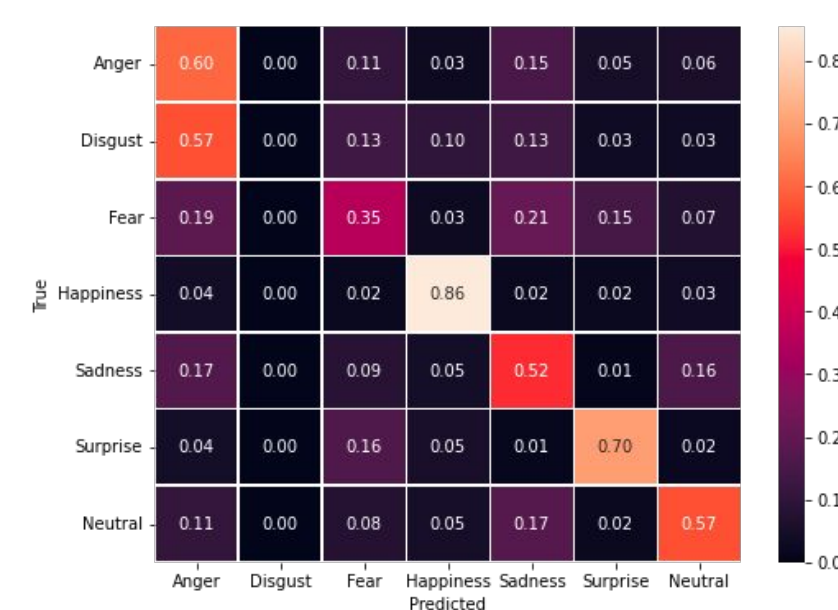
### RESULTS



Normalized confusion matrix for FERNet



Normalized confusion matrix for ResNet



Normalized confusion matrix for ViT

	AlexNet (FERNet)	ResNet	ViT
Training time	~13 hours	~5 hours	37 minutes (smaller dataset)
Number of parameters	29,069,064	223,847	86,394,631
Final test accuracy (Human: 65±5%)	50.4%	60%	61%
Final training accuracy	72.4%	63.8%	N/A
Number of epochs	30	30	6

### CONCLUSION

Both ResNet and ViT showed significantly higher test accuracy scores when compared with the older AlexNet model. These models, despite their implementation differences, both performed relatively well with test accuracy scores in the range of the human accuracy of 65.5+5% on the FER2013 dataset. Notably, ResNet was able to correctly predict the "disgust" label which is impressive because of its substantially smaller sample size compared to other labels. On the other hand, the ViT model was able to achieve the highest test accuracy score despite being trained on a smaller training dataset.

### REFERENCES

arXiv:1307.0414  
arXiv:1512.03385  
arXiv:2010.11929